# Future of the Earth
# & its Decisive Factors

Columbus Foust

First Edition, 2012

# Table of Contents

# Chapter-1

# Introduction to Future of the Earth



Conjectured illustration of the scorched Earth after the Sun has entered the red giant phase, seven billion years from now.

The **future of the Earth** will be determined by a variety of factors, including increases in the luminosity of the Sun, loss of heat energy from the Earth's core, perturbations by the other bodies in the Solar System and the biochemistry at the Earth's surface. Milankovitch theory predicts the planet will continue to undergo glaciation cycles because of eccentricity, axial tilt, and precession of the Earth's orbit. As part of the ongoing supercontinent cycle, plate tectonics will probably result in a supercontinent in 250 million–350 million years. Some time in the next 1.5 billion–4.5 billion years, the axial tilt of the Earth may begin to undergo chaotic variations, with changes in the axial tilt of up to 90°.

One billion to two billion years in the future, the steady increase in solar radiation caused by the helium build-up at the core of the Sun will result in the loss of the oceans and the cessation of continental drift. Four billion years from now, the increase in the Earth's surface temperature will cause a runaway greenhouse effect. By that point, most if not all the life on the surface will be extinct. The most likely ultimate fate of the planet is absorption by the Sun in about 7.5 billion years, after the star has entered the red giant phase and expanded to cross the planet's orbit.

# Human influence

Humans now play a key role in the biosphere, with the large human population dominating many of Earth's ecosystems. This has resulted in a widespread, ongoing extinction of other species during the present geological epoch, now known as the Holocene extinction. The large scale loss of species caused by human influence since the 1950s has been called a biotic crisis, with an estimated 10% of the total species lost as of 2007. At current rates, about 30% of species are at risk of extinction in the next hundred years. The Holocene extinction event is the result of habitat destruction, the widespread distribution of invasive species, and hunting and climate change. In the present day, human activity has had a significant impact on the surface of the planet. More than a third of the land surface has been modified by human actions, and humans use about 20% of global primary production. The concentration of carbon dioxide in the atmosphere has increased by close to 30% since the start of the Industrial Revolution.

The consequences of a persistent biotic crisis have been predicted to last for at least five million years. It could result in a decline in biodiversity and homogenization of biotas, accompanied by a proliferation of species that are opportunistic, such as pests and weeds. Novel species may also emerge; in particular taxa that prosper in human-dominated ecosystems may rapidly diversify into many new species. Microbes are likely to benefit from the increase in nutrient-enriched environmental niches. However, no new species of existing large vertebrates are likely to arise and food chains will probably be shorter.

# Orbit and rotation

The gravitational perturbations of the other planets in the Solar System combine to modify the orbit of the Earth and the orientation of its spin axis. These changes can influence the planetary climate.

## Glaciation

Historically, there have been cyclical periods of glaciation in which ice sheets covered the higher latitudes of the continents. The Milankovitch theory predicts that glaciation occurs because of astronomical factors in combination with climate feedback mechanisms and plate tectonics. The primary astronomical drivers are a higher than normal orbital eccentricity, a low axial tilt (or obliquity), and the alignment of summer solstice with the aphelion. Each of these effects occur cyclically. For example, the

eccentricity changes over time cycles of about 100,000 and 400,000 years, with the value ranging from less than 0.01 up to 0.05. This is equivalent to a change of the semiminor axis of the planet's orbit from 99.95% of the semimajor axis to 99.88%, respectively.

At present the Earth is in an interglacial period, which would normally be expected to end in about 25,000 years. The current rate of increased carbon dioxide release into the atmosphere by humans may delay the onset of the next period of glaciation until at least 50,000–130,000 years from now. However, a global warming period of finite duration (based on the assumption that fossil fuel use will cease by the year 2200) will probably only impact the glaciation cycle for about 5,000 years. Thus, a brief period of global warming induced through a few centuries worth of greenhouse gas emission would only have a limited impact in the long term.

## Obliquity



The rotational offset of the tidal bulge exerts a net torque on the Moon, boosting it while slowing the Earth's rotation.

The tidal acceleration of the Moon slows the rotation rate of the Earth and increases the Earth-Moon distance. Other effects that can dissipate the Earth's rotational energy are friction between the core and mantle, tides in the atmosphere, convection in the mantle, and climate changes that can increase or decrease the ice load at the poles. These combined effects are expected to increase the length of the day by more than 1.5 hours over the next 250 million years, and to increase the obliquity by about a half degree. The distance to the Moon will increase by about 1.5 Earth radii during the same period.

Based on computer models, the presence of the Moon appears to stabilize the obliquity of the Earth, which may help the planet to avoid dramatic climate changes. This stability is achieved because the Moon increases the precession rate of the Earth's spin axis, thereby avoiding resonances between the precession of the spin and precession frequencies of the

ascending node of the planet's orbit. (That is, the precession motion of the ecliptic.) However, as the semimajor axis of the Moon's orbit continues to increase in the future, this stabilizing effect will diminish. At some point perturbation effects will probably cause chaotic variations in the obliquity of the Earth, and the axial tilt may change by angles as high as 90° from the plane of the orbit. This is expected to occur within about 1.5–4.5 billion years, although the exact time is unknown.

A high obliquity would probably result in dramatic changes in the climate and may destroy the planet's habitability. When the axial tilt of the Earth reaches 54°, the equator will receive less radiation from the Sun than the poles. The planet could remain at an obiliquity of 60° to 90° for periods as long as 10 million years.

# Plate tectonics

Pangaea was the last supercontinent to form before the present.

The theory of plate tectonics demonstrates that the continents of the Earth are moving across the surface at the rate of a few centimeters per year. This is expected to continue, causing the plates to relocate and collide. Continental drift is facilitated by two factors: the energy generation within the planet and the presence of a hydrosphere. With the loss of either of these, continental drift will come to a halt. The production of heat through radiogenic processes is sufficient to maintain mantle convection and plate subduction for at least the next 1.1 billion years.

At present, the continents of North and South America are moving westward from Africa and Europe. Researchers have produced several scenarios about how this will continue in the future. These geodynamic models can be distinguished by the subduction flux, whereby the oceanic crust moves under a continent. In the introversion model, the younger, interior, Atlantic ocean becomes preferentially subducted and the current migration of North and South America is reversed. In the extroversion model, the older, exterior, Pacific ocean remains preferentially subducted and North and South America migrate toward eastern Asia.

As the understanding of geodynamics improves, these models will be subject to revision. In 2008, for example, a computer simulation was used to predict that a reorganization of the mantle convection will occur, causing a supercontinent to form around Antarctica.

Regardless of the outcome of the continental migration, the continued subduction process causes water to be transported to the mantle. After a billion years from the present, a geophysical model gives an estimate that 27% of the current ocean mass will have been subducted. If this process were to continue unmodified into the future, the subduction and release would reach a point of stability after 65% of the current ocean mass has been subducted.

## Introversion

Christopher Scotese and his colleagues have mapped out the predicted motions several hundred million years into the future as part of the Paleomap Project. In their scenario, 50 million years from now the Mediterranean sea may vanish and the collision between Europe and Africa will create a long mountain range extending to the current location of the Persian Gulf. Australia will merge with Indonesia, and Baja California will slide northward along the coast. New subduction zones may appear off the eastern coast of North and South America, and mountain chains will form along those coastlines. To the south, the migration of Antarctica to the north will cause all of its ice sheets to melt. This, along with the melting of the Greenland ice sheets, will raise the average ocean level by 90 metres (300 ft). The inland flooding of the continents will result in climate changes.

As this scenario continues, by 100 million years from the present the continental spreading will have reached its maximum extent and the continents will then begin to coalesce. In 250 million years, North America will collide with Africa while South America will wrap around the southern tip of Africa. The result will be the formation of a new supercontinent (sometimes called Pangaea Ultima), with the Pacific Ocean stretching across half the planet. The continent of Antarctica will reverse direction and return to the South Pole, building up a new ice cap.

## Extroversion

The first scientist to extrapolate the current motions of the continents was Canadian geologist Paul F. Hoffman of Harvard University. In 1992, Hoffman predicted that continents of North and South America would continue to advance across the Pacific

Ocean, pivoting about Siberia until they begin to merge with Asia. He dubbed the resulting supercontinent, Amasia. Later, in the 1990s, Roy Livermore calculated a similar scenario. He predicted that Antarctica would start to migrate northward, and east Africa and Madagascar would move across the Indian Ocean to collide with Asia.

In an extroversion model, the closure of the Pacific Ocean would be complete by about 350 million years. This marks the completion of the current supercontinent cycle, wherein the continents split apart and then rejoin each other about every 400–500 million years. Once the supercontinent is built, plate tectonics may enter a period of inactivity as the rate of subduction drops by an order of magnitude. This period of stability could cause an increase in the mantle temperature at the rate of 30–100 K every 100 million years, which is the minimum lifetime of past supercontinents. As a consequence, volcanic activity may increase.

### Supercontinent

The formation of a supercontinent can dramatically affect the environment. The collision of plates will result in mountain building, thereby shifting weather patterns. Sea levels may fall because of increased glaciation. The rate of surface weathering can rise, resulting in an increase in the rate that organic material is buried. Supercontinents can cause a drop in global temperatures and an increase in atmospheric oxygen. These changes can result in more rapid biological evolution as new niches emerge. This, in turn, can affect the climate, further lowering temperatures.

The formation of a supercontinent insulates the mantle. The flow of heat will be concentrated, resulting in volcanism and the flooding of large areas with basalt. Rifts will form and the supercontinent will split up once more.

# Solar evolution

The energy generation of the Sun is based upon thermonuclear fusion of hydrogen into helium. This occurs in the core region of the star using the proton–proton chain reaction process. Because there is no convection in the solar core, the fusion process results in a steady build-up of helium. The temperature at the core of the Sun is too low for nuclear fusion of helium atoms through the triple-alpha process, so these atoms do not contribute to the net energy generation that is needed to maintain hydrostatic equilibrium of the Sun.

At present, nearly half the hydrogen at the core has been consumed, with the remainder consisting primarily of helium. To compensate for the steadily decreasing number of hydrogen atoms per unit mass, the core temperature of the Sun has gradually increased through a rise in pressure. This has caused the remaining hydrogen to undergo fusion at a more rapid rate, thereby generating the energy needed to maintain the equilibrium. The result has been a steady increase in the energy output of the Sun. This increase can be approximated by the formula:

$$L(t) = \left[1 + \frac{2}{5}\left(1 - \frac{t}{t_{Sun}}\right)\right]^{-1} L_{Sun}$$

where $t$ is a time period less than or equal to the present time $t_{Sun}$, $L(t)$ is the luminosity at time $t$, and $L_{Sun}$ is the current solar luminosity.

When the Sun first became a main sequence star, it radiated only 70% of the current luminosity. The luminosity has increased in a nearly linear fashion to the present, increasing by 1% every 110 million years. Likewise, in three billion years the Sun is expected to be 33% more luminous. The hydrogen fuel at the core will finally be exhausted in 4.8 billion years, when the Sun will be 67% more luminous than at present. Thereafter the Sun will continue to burn hydrogen in a shell surrounding its core, until the increase in luminosity reaches 121% of the present value. This marks the end of the Sun's main sequence lifetime, and thereafter it will evolve into a red giant.
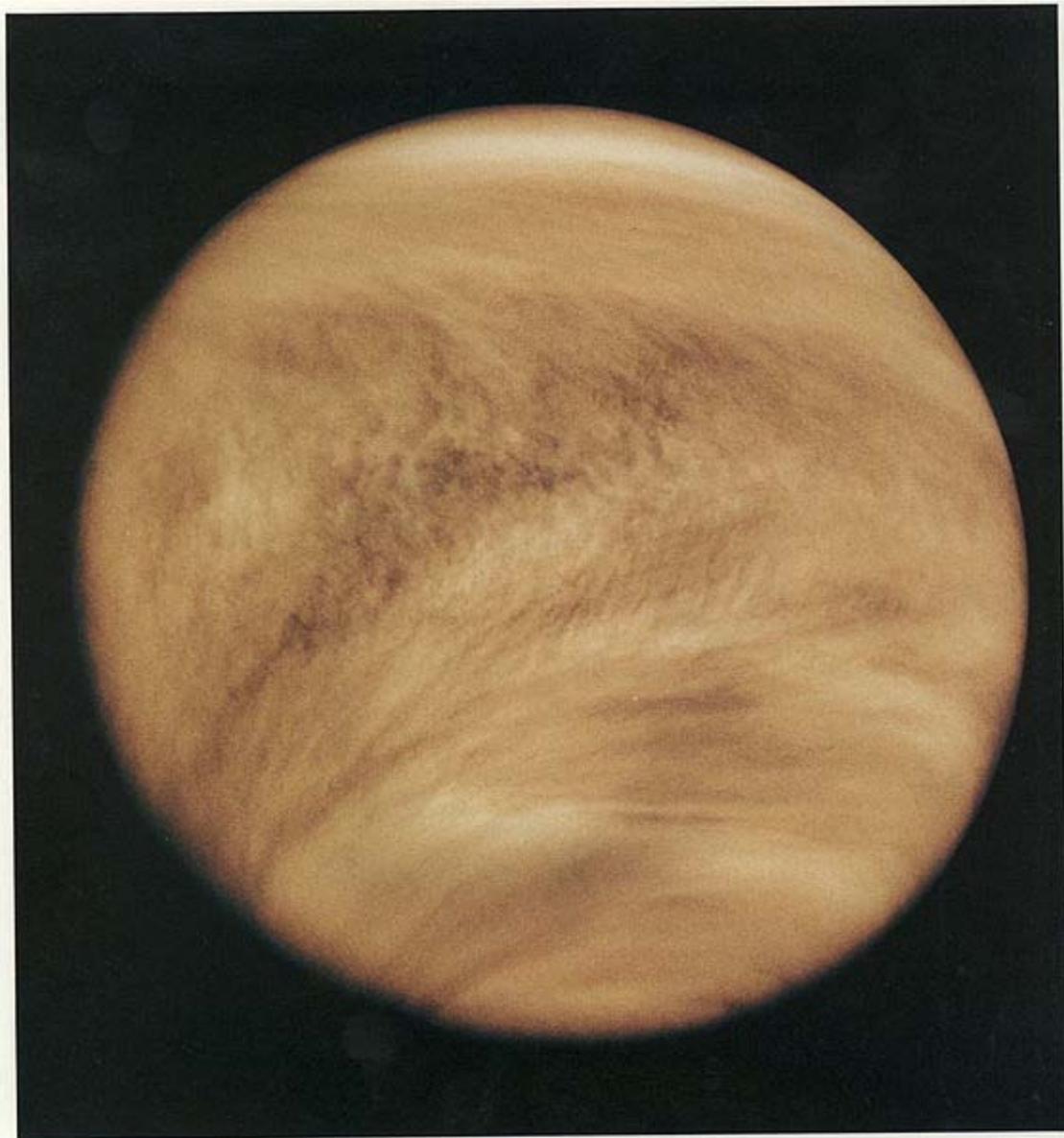
## Climate impact

As the global temperature of the Earth climbs because of the rising luminosity of the Sun, the rate of weathering of silicate minerals will increase. This in turn will decrease the level of carbon dioxide in the atmosphere. Within the next 600 million years from the present, the concentration of $CO_2$ will fall below the critical threshold needed to sustain $C_3$ photosynthesis: about 50 parts per million. At this point, trees and forests in their current forms will no longer be able to survive. However, $C_4$ carbon fixation can continue at much lower concentrations, down to above 10 parts per million. Thus plants using $C_4$ photosynthesis may be able to survive for at least 0.8 billion years and possibly as long as 1.2 billion years from now, after which rising temperatures will make the biosphere unsustainable. Currently, $C_4$ plants represent about 5% of Earth's plant biomass and 1% of its known plant species. For example, about 50% of all grass species (Poaceae) use the $C_4$ photosynthetic pathway, as do many species in the herbaceous family Amaranthaceae.

When the levels of carbon dioxide fall to the limit where photosynthesis is barely sustainable, the proportion of carbon dioxide in the atmosphere is expected to oscillate up and down. This will allow land vegetation to reappear each time the level of carbon dioxide rises due to tectonic activity and animal life. However, the long term trend is for the plant life on land to die off altogether as most of the remaining carbon in the atmosphere becomes sequestered in the Earth. Some microbes are capable of photosynthesis at concentrations of $CO_2$ of a few parts per million, so these life forms would probably disappear only because of rising temperatures and the loss of the biosphere.

In their work *The Life and Death of Planet Earth*, authors Peter D. Ward and Donald Brownlee have argued that some form of animal life may continue even after most of the Earth's plant life has disappeared. Initially, they expect that some insects, lizards, birds and small mammals may persist, along with sea life. Without oxygen replenishment by plant life, however, they believe that the animals would probably die off from asphyxiation within a few million years. Even if sufficient oxygen were to remain in the

atmosphere through the persistence of some form of photosynthesis, the steady rise in global temperature would result in a gradual loss of biodiversity. Much of the surface would become a barren desert and life would primarily be found in the oceans.

Once the solar luminosity is 10% higher than its current value, the average global surface temperature reaches 320 K (47 °C). The atmosphere will become a humid greenhouse leading to a runaway evaporation of the oceans. At this point, models of the Earth's future environment demonstrate that the stratosphere would contain increasing levels of water. These water molecules will be broken down through photodissociation by solar ultraviolet radiation, allowing hydrogen to escape the atmosphere. The net result would be a loss of the world's sea water in about 1.1 billion years from the present.
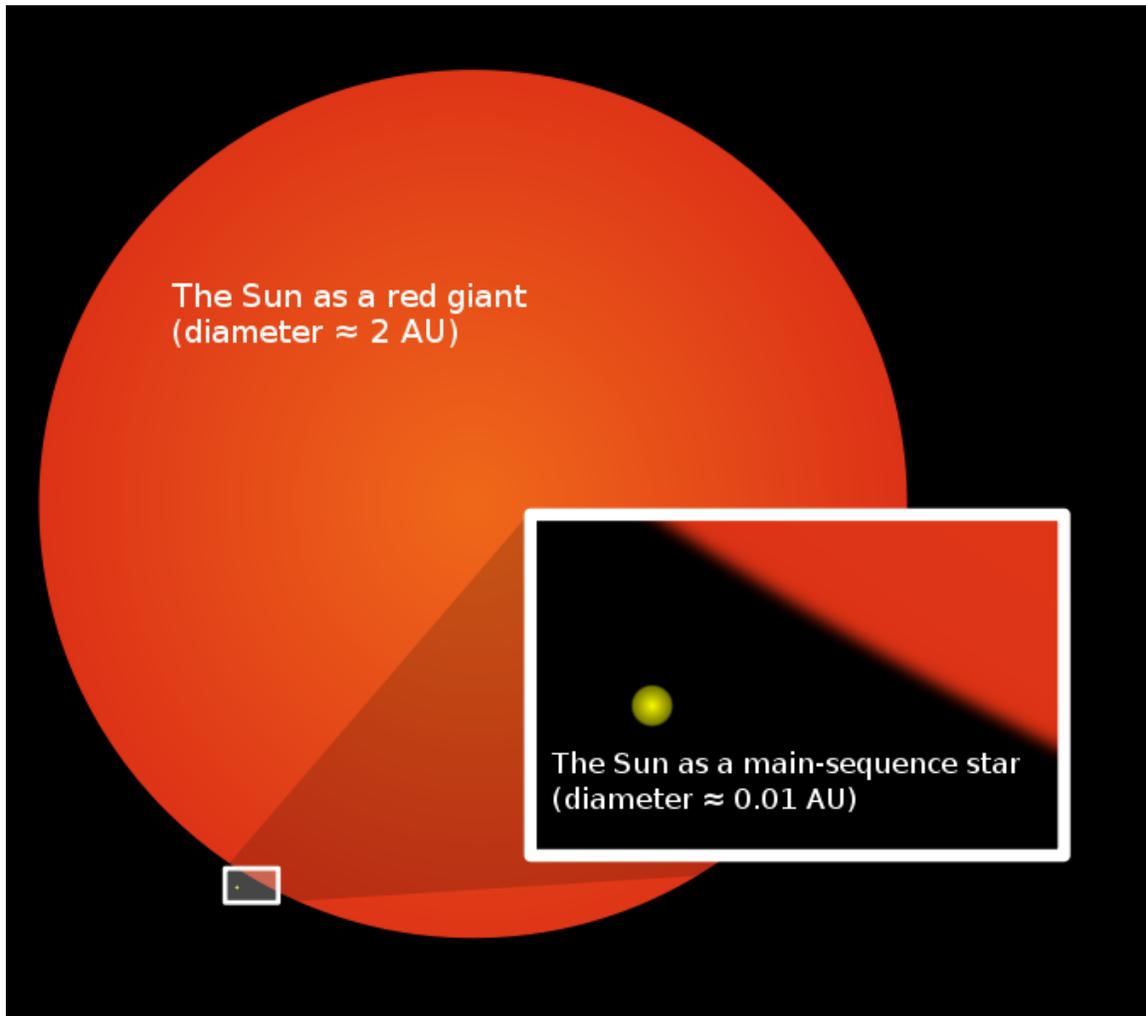


The atmosphere of Venus is in a "supergreenhouse" state.

Still, there will continue to be some reservoirs at the surface as water is steadily released from the deep crust and mantle. Some water may be retained at the poles and there may be occasional rainstorms, but for the most part the planet would be a dry desert. What happens next depends on the level of tectonic activity. The release of carbon dioxide by volcanic eruption may eventually cause the atmosphere to enter a "supergreenhouse" state like that of the planet Venus. However, without surface water, plate tectonics would probably come to a halt and most of the carbonates would remain securely buried.

The loss of the oceans could be delayed until two billion years in the future if the total atmospheric pressure were to decline. A lower atmospheric pressure would reduce the greenhouse effect, thereby lowering the surface temperature. This could occur if natural processes were to remove the nitrogen from the atmosphere. Studies of organic sediments has shown that at least 100 kilopascals (1 bar) of nitrogen has been removed from the atmosphere over the past four billion years; enough to effectively double the current atmospheric pressure if it were to be released. This rate of removal would be sufficient to counter the effects of increasing solar luminosity for the next two billion years. However, beyond that point, the amount of water in the lower atmosphere will have risen to 40% and the runaway moist greenhouse will commence.

A runaway greenhouse effect will take place when the luminosity from the Sun reaches 40% more than its current value, four billion years from now. The atmosphere will heat up and the surface temperature will rise. However, most of the atmosphere will be retained until the Sun has entered the red giant stage.

**Red giant stage**



The size of the current Sun (now in the main sequence) compared to its estimated size during its red giant phase.

Once the Sun changes from burning hydrogen at the core to burning hydrogen around a shell, the core will start to contract and the outer envelope will expand. The total luminosity will steadily increase over the next billion years until it reaches 2,730 times the Sun's current luminosity at the age of 12.167 billion years. During this phase the Sun will undergo mass loss, with about 33% of its total mass shed with the solar wind. The loss of mass will mean that the orbits of the planets will expand. The orbital distance of the Earth will increase to at most 150% of its current value.

The most rapid part of the Sun's expansion into a red giant occurs during the final stages, when the Sun is about 12 billion years old. It is likely to expand to swallow both Mercury and Venus, reaching a maximum radius of 1.2 astronomical units (180 Gm). The Earth will interact tidally with the Sun's outer atmosphere, which would serve to decrease the

orbital radius. Drag from the chromosphere of the Sun would also reduce the Earth's orbit. These effects will act to counterbalance the mass loss by the Sun, and the Earth will most likely be engulfed by the sun.

By the time the Sun begins to grow as a red giant, the orbit of the Moon will have expanded until it takes 47 days to complete. The drag from the solar atmosphere may cause the orbit of the Moon to decay. Once the orbit of the Moon closes to a distance of 18,470 km, it will cross the Earth's Roche limit. The tidal interaction with the Earth will break apart the Moon, turning it into a ring system. Most of the orbiting ring will then begin to decay, and the debris will impact the Earth. Hence, even if the Earth is not swallowed up by the Sun, the planet may be left moonless.

# Chapter-2

# Conservation Biology



Efforts are being taken to preserve the natural characteristics of Hopetoun Falls, Australia while continuing to allow visitor access

**Conservation biology** is the scientific study of the nature and status of Earth's biodiversity with the aim of protecting species, their habitats, and ecosystems from excessive rates of extinction. It is an interdisciplinary subject drawing on sciences, economics, and the practice of natural resource management.

## History of term

The term *conservation biology* was introduced as the title of a conference held at the University of California in La Jolla, California in 1978 organized by biologists Bruce Wilcox and Michael Soulé. The meeting was prompted by the concern among scientists

over tropical deforestation, disappearing species, eroding genetic diversity within species. The conference and proceedings that resulted sought to bridge a gap existing at the time between theory in ecology and population biology on the one hand and conservation policy and practice on the other. Conservation biology and the concept of biological diversity (biodiversity) emerged together, helping crystallize the modern era of conservation science and policy.

# Description

The rapid decline of established biological systems around the world means that conservation biology is often referred to as a "Discipline with a deadline". Conservation biology is tied closely to ecology in researching the dispersal, migration, demographics, effective population size, inbreeding depression, and minimum population viability of rare or endangered species. Conservation biology is concerned with phenomena that affect the maintenance, loss, and restoration of biodiversity and the science of sustaining evolutionary processes that engender genetic, population, species, and ecosystem diversity. The concern stems from estimates suggesting that up to 50% of all species on the planet will disappear within the next 50 years, which has contributed to poverty, starvation, and will reset the course of evolution on this planet.
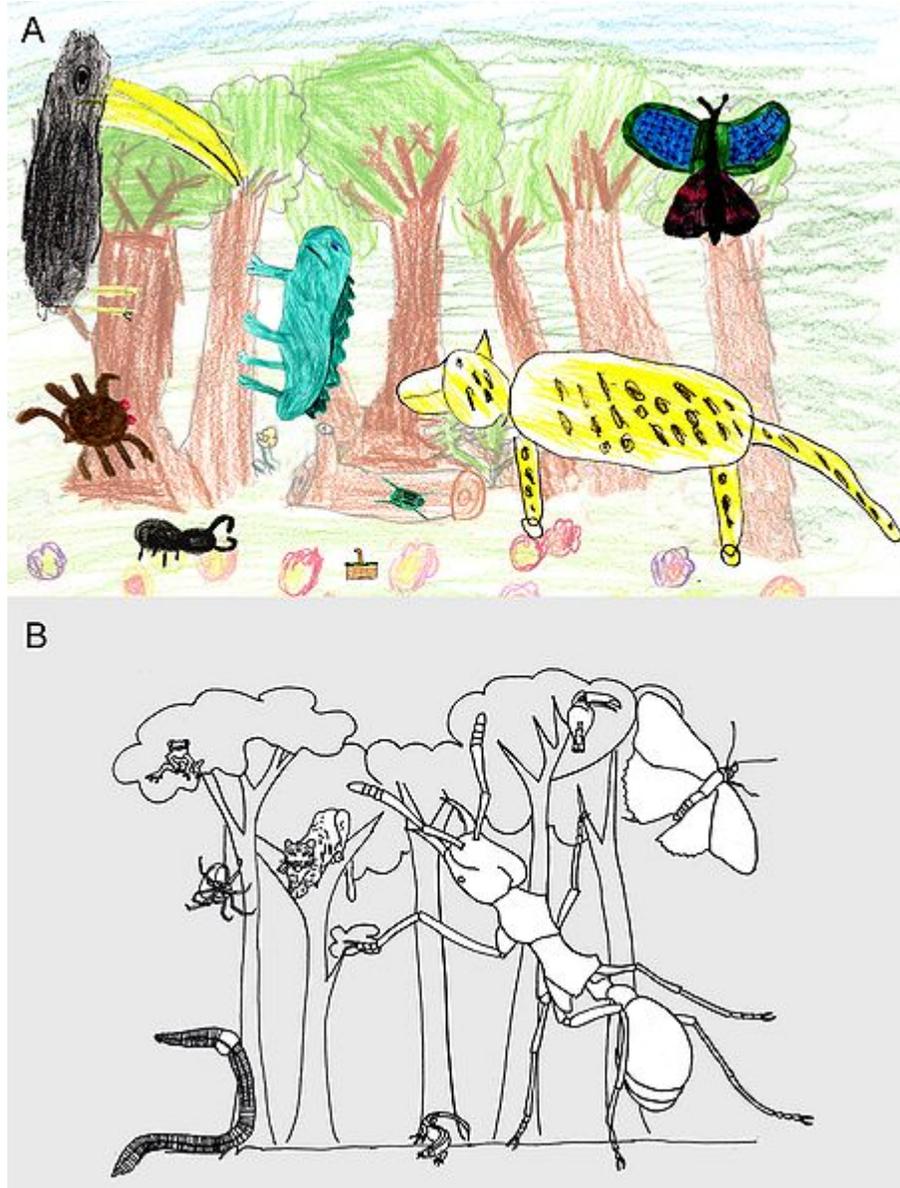
Conservation biologists research and educate on the trends and process of biodiversity loss, species extinctions, and the negative affect this is having on our capabilities to sustain the well-being of human society. Conservation biologists work in the field and office, in government, universities, non-profit organizations and industry. They are funded to research, monitor, and catalog every angle of the earth and its relation to society. The topics are diverse, because this is an interdisciplinary network with professional alliances in the biological as well as social sciences. Those dedicated to the cause and profession advocate for a global response to the current biodiversity crisis based on morals, ethics, and scientific reason. Organizations and citizens are responding to the biodiversity crisis through conservation action plans that direct research, monitoring, and education programs that engage concerns at local through global scales.

# Context and trends

Conservation biologists study trends and process from the paleontological past to the ecological present as they gain an understanding of the context related to species extinction. It is generally accepted that there have been five major global mass extinctions that register in Earth's history. These include: the Ordovician (440 mya), Devonian (370 mya), Permian–Triassic (245 mya), Triassic–Jurassic (200 mya), and Cretaceous (65 mya) extinction spasms. Within the last 10,000 years, human influence over the Earth's ecosystems has been so extensive that scientists have difficulty estimating the number of species lost; that is to say the rates of deforestation, reef destruction, wetland draining and other human acts are proceeding much faster than human assessment of species. The latest *Living Planet Report* by the World Wide Fund

for Nature estimates that we have exceeded the bio-regenerative capacity of the planet, requiring 1.5 Earths to support the demands placed on our natural resources.

## Sixth extinction



An art scape image showing the relative importance of animals in a rain forest through a summary of (a) child's perception compared with (b) a scientific estimate of the importance. The size of the animal represents its importance. The child's mental image places importance on big cats, birds, butterflies, and then reptiles versus the actual dominance of social insects (such as ants).

Conservation biologists are dealing with and have published evidence from all corners of the planet indicating that humanity may be living the sixth and greatest planetary

extinction event. It has been suggested that we are living in an era of unprecedented numbers of species extinctions, also known as the Holocene extinction event. The global extinction rate may be approximately 100,000 times higher than the natural background extinction rate. It is estimated that two-thirds of all mammal genera and one-half of all mammal species weighing at least 44 kilograms (97 lb) have gone extinct in the last 50,000 years. It is speculated that this sixth extinction period is unique because it would be the first major extinction to be caused by another biotic agent over the course of the Earth's 4 billion year history. The Global Amphibian Assessment reports that amphibians are declining on a global scale faster than any other vertebrate group, with over 32% of all surviving species being threatened with extinction. The surviving populations are in continual decline in 43% of those that are threatened. Since the mid-1980s the actual rates of extinction have exceeded 211 times rates measured from the fossil record. However, "The current amphibian extinction rate may range from 25,039 to 45,474 times the background extinction rate for amphibians." The global extinction trend occurs in every major vertebrate group that is being monitored. For example, 23% of all mammals and 12% of all birds are Red Listed by the International Union for Conservation of Nature (IUCN), meaning they too are threatened with extinction.

## Status of oceans and reefs

Global assessments of coral reefs of the world continue to report drastic and rapid rates of decline. By 2000, 27% of the world's coral reef ecosystems had effectively collapsed. The largest period of decline occurred in a dramatic "bleaching" event in 1998, where approximately 16% of all the coral reefs in the world disappeared in less than a year. *Coral bleaching* is caused by a mixture of environmental stresses, including increases in ocean temperatures and acidity, causing both the release of symbiotic algae and death of corals. Decline and extinction risk in coral reef biodiversity has risen dramatically in the past ten years. The loss of coral reefs, which are predicted to go extinct in the next century, will have huge economic impacts, threatens the balance of global biodiversity, and endangers food security for hundreds of millions of people. Conservation biology plays an important role in international agreements covering the world's oceans (and other issues pertaining to biodiversity, e.g.).

**These predictions will undoubtedly appear extreme, but it is difficult to imagine how such changes will not come to pass without fundamental changes in human behavior.**

J.B. Jackson

The oceans are threatened by acidification due to an increase in $CO_2$ levels. This is a most serious threat to societies relying heavily upon oceanic natural resources. A concern is that the majority of all marine species will not be able to evolve or acclimate in response to the changes in the ocean chemistry.

The prospects of averting mass extinction seems unlikely when "[...] 90% of all of the large (average approximately $\geq$50 kg), open ocean tuna, billfishes, and sharks in the ocean" are reportedly gone. Given the scientific review of current trends, the ocean is

predicted to have few surviving multi-cellular organisms with only microbes left to dominate marine ecosystems.

## Insects and other groups

There are serious concerns also being hailed from taxonomic groups that do not receive the same degree of social attention or attract funds as the vertebrates do, including fungi, lichen, plant and insect communities where the vast majority of biodiversity is represented. Insect conservation, in particular, is of pivotal importance for conservation biology. The value of insects in the biosphere is enormous because they outnumber all other living groups in measure of species richness. The greatest bulk of biomass on land is found in plants, which is sustained by insect relations. This great ecological value of insects is countered by a society that oftentimes reacts negatively toward these aesthetically 'unpleasant' creatures.

One area of concern in the insect world that has caught the public eye is the mysterious case of missing honey bees (*Apis mellifera*). Honey bees provide an indispensable ecological services through their acts of pollination supporting a huge variety of agriculture crops. The sudden disappearance of bees leaving empty hives or colony collapse disorder (CCD) is not uncommon. However, in 16-month period from 2006 through 2007, 29% of 577 beekeepers across the United States reported CCD losses in up to 76% of their colonies. This sudden demographic loss in bee numbers is placing a strain on the agricultural sector. The cause behind the massive declines is puzzling scientists. Pests, pesticides, and global warming are all being considered as possible causes.

Another highlight that links conservation biology to insects, forests, and climate change is the mountain pine beetle (*Dendroctonus ponderosae*) epidemic of British Columbia, Canada, which has infested 470,000 km$^2$ (180,000 sq mi) of forested land since 1999. An action plan has been prepared by the Government of British Columbia to address this problem.

This impact [*pine beetle epidemic*] converted the forest from a small net carbon sink to a large net carbon source both during and immediately after the outbreak. In the worst year, the impacts resulting from the beetle outbreak in British Columbia were equivalent to 75% of the average annual direct forest fire emissions from all of Canada during 1959–1999.
—Kurz *et al*.

# Conservation biology of parasites



The capture, captive breeding, and reintroduction of California Condor into the wild was the most expensive species conservation project in United States history. The bird was saved from extinction while its louse *Colpocephalum californici* went extinct.

A large proportion of living species on Earth live a parasitic way of life. Parasites have traditionally been seen as targets of eradication efforts, and they have often been overlooked in conservation efforts. In the case of parasites living in the wild – and thus harmless to humans and domesticated animals – this view is changing.

# Endangered parasite species

A note published in 1990 pointed out that the captive breeding and reintroduction program to save the black-footed ferret would cause the loss of its specific parasites and demanded *"equal rights for parasites!"*. Then a paper in 1992 has warned that not only the loss of certain host species from the wild, but even host population bottlenecks or the fragmentation of host populations would predictably lead to the extinction of several host specific parasite species. It also noted that parasites are not only components of biodiversity by definition, but they also exert selective pressures upon their host populations that increase host genetic diversity. Firstly, this view met with open scepticism. Soon after, it became clear that the co-extinction of hosts and their specific parasites is likely to increase the current estimates of extinction rates significantly. A decade later, a study focusing on some highly host-specific groups (such as fig wasps, parasites, butterflies, and myrmecophil butterflies) estimated the number of co-endangered species (i.e. endagered by the endagered status of the host) at about 6300. Other authors argued that host specific parasite faunae have an unexpected advantage for conservation scientists. Their genealogies and population genetic patterns may help to illuminate their hosts' evolutionary and demographic history. Recently, scientists suggested that rich parasite faunae are inevitably needed for healthy ecosystem functioning and also that parasites and mutualists are the most endangered species on Earth. Even veterinarians have started to argue about the conservational values of parasite species.

# Example: extinct avian lice

The list below follows that of Mey (2005)

- *Acutifrons caracarensis* parasite of the extinct Guadalupe Caracara *(Caracara lutosa)*, Guadalupe Island, Mexico;
- *Longimenopon dominicanum* parasite of the extinct Guadalupe Storm-petrel, *Oceanodroma macrodactyla*, Guadalupe Island, Mexico;
- *Campanulotes defectus* parasite of the extinct Passenger Pigeon *(Ectopistes migratorius)*, North-America;
- (*Columbicola extinctus* another parasite of the extinct Passenger Pigeon (*Ectopistes migratorius*). Interestingly, recent taxonomic studies show that it was conspecific with the lice living on Band-tailed Pigeon *(Columba fasciata)*, thus it is not extinct as a species) ;
- *Rallicola piageti* parasite of the extinct New Caledonian Rail *(Gallirallus lafresnayanus)*, New-Caledonia;
- *Halipeurus raphanus* parasite of the extinct Guadalupe Storm-petrel *(Oceanodroma macrodactyla)*, Guadelupe Island, Mexico;
- *Puffinoecus jamaicensis* parasite of the extinct Jamaica Petrel *(Pterodroma caribbaea)*, Jamaica;
- *Nitzschiella hemiphagae* parasite of the extinct Norfolk Island Pigeon (*Hemiphaga novaeseelandiae spadicea*), Norfolk Island, New-Zealand;

- *Patellinirmus restinctus* parasite of the extinct Norfolk Island Pigeon (*Hemiphaga novaeseelandiae spadicea*), Norfolk Island, New-Zealand;
- *Rallicola extinctus* parasite of the extinct Huia *(Heteralocha acutirostris)*, New-Zealand;
- *Philopteroides xenicus* parasite of the extinct Bushwren *(Xenicus longipes)*, New-Zealand;
- *Psittacobrosus bechsteini* parasite of the extinct Cuban Red Macaw *(Ara tricolor)*, Cuba;
- *Colpocephalum californici*, parasite of the California Condor *(Gymnogyps californianus)*. The host have been saved by captive breeding and repatriation programs, however, the parasite have been lost, either spontaneously or perhaps exterminated by wildlife vets.



Larvae of the Guinea Worm: probably the next species to exterminate.

# Extermination by purpose

Naturally, medical (and veterinary) science and practice aim to exterminate parasites and pathogens living in humans (and in domesticated animals). In case of the few highly host-specific pathogens, this equals the extinction of the pathogen species. Throughout human history, however, only a single one species, i.e. smallpox virus, was eradicated from the Globe. The last cases of smallpox occurred 1978. However, secured stocks still exist in the United States and Russia for defensive purposes such as developing new vaccines,

antiviral drugs, and diagnostic tests. It is not known whether or not these superpowers have shared their stocks with some of their allies during the Cold War.

A second candidate for purposeful extermination is the Guinea Worm *(Dracunculus medinensis)*. Once widespread across some 20 nations of Africa and Asia, the parasite nowadays is much withdrown occurring only in a few countries of Sub-Saharan Africa. Prevalent civil wars in the region, such as the War in Darfur have ensured the survival of this species up to the present.

## Threats to biodiversity

Many of the threats to biodiversity, including disease and climate change, are reaching inside borders of protected areas, leaving them 'not-so protected' (e.g. Yellowstone National Park). Climate change, for example, is often cited as a serious threat in this regard, because there is a feedback loop between species extinction and the release of carbon dioxide into the atmosphere. Ecosystems store and cycle large amounts of carbon to regulate global conditions. The effects of global warming adds a catastrophic threat toward a mass extinction of global biological diversity. The extinction threat is estimated to range from 15 to 37 percent of all species by 2050, or 50 percent of all species over the next 50 years.

Some of the most significant and insidious threats to biodiversity and ecosystem processes include climate change, mass agriculture, deforestation, overgrazing, slash-and-burn agriculture, urban development, wildlife trade, light pollution and pesticide use. Habitat fragmentation poses one of the more difficult challenges, because the global network of protected areas only covers 11.5% of the Earth's surface. A significant consequence of fragmentation and lack of linked protected areas is the reduction of animal migration on a global scale. Considering that billions of tonnes of biomass are responsible for nutrient cycling across the earth, the reduction of migration is a serious matter for conservation biology.
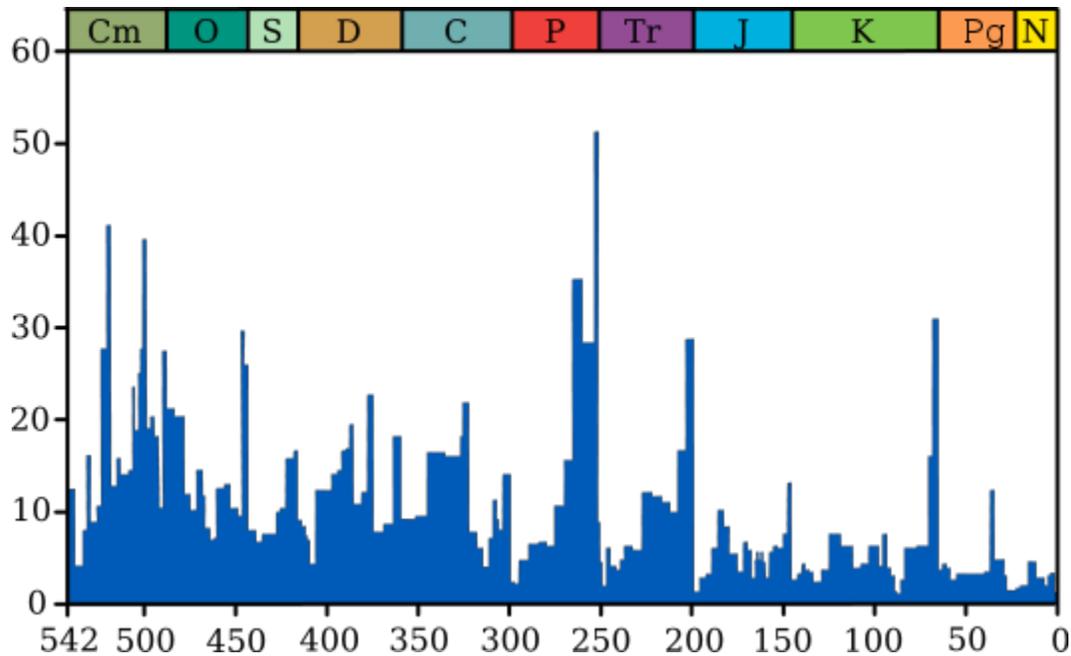
Human activities are associated directly or indirectly with nearly every aspect of the current extinction spasm.

Wake and Vredenburg

These figures do not imply, however, that human activities must necessarily cause irreparable harm to the biosphere. With conservation management and planning for biodiversity at all levels, from genes to ecosystems, there are examples where humans mutually coexist in a sustainable way with nature. However, it may be too late for human intervention to reverse the current mass extinction.

# Concepts and foundations

## Measuring extinction rates



The five major extinction spasms measured by extinction levels in marine animal genera through time. Blue graph shows apparent percentage (not absolute number) of extinctions during any given time interval.

Extinction rates are measured in a variety of ways. Conservation biologists measure and apply statistical measures of fossil records, rates of habitat loss, and a multitude of other variables such as loss of biodiversity as a function of the rate of habitat loss and site occupancy to obtain such estimates. The Theory of Island Biogeography is possibly the most significant contribution toward the scientific understanding of both the process and how to measure the rate of species extinction. The current background extinction rate is estimated to be one species every few years.

The measure of ongoing species loss is made more complex by the fact that most of the Earth's species have not been described or evaluated. Estimates vary greatly on how many species actually exist (estimated range: 3,600,000-111,700,000) to how many have received a species binomial (estimated range: 1.5-8 million). Less than 1% of all species that have been described have been studied beyond simply noting its existence. From these figures, the IUCN reports that 23% of vertebrates, 5% of invertebrates and 70% of plants that have been evaluated are designated as endangered or threatened.

## Systematic conservation planning

Systematic conservation planning is an effective way to seek and identify efficient and effective types of reserve design to capture or sustain the highest priority biodiversity values and to work with communities in support of local ecosystems. Margules and Pressey identify six interlinked stages in the systematic planning approach:

1. Compile data on the biodiversity of the planning region
2. Identify conservation goals for the planning region
3. Review existing conservation areas
4. Select additional conservation areas
5. Implement conservation actions
6. Maintain the required values of conservation areas

Conservation biologists regularly prepare detailed conservation plans for grant proposals or to effectively coordinate their plan of action and to identify best management practices (e.g.). Systematic strategies generally employ the services of Geographic Information Systems to assist in the decision making process.

## Conservation biology as a profession

The Society for Conservation Biology is a global community of conservation professionals dedicated to advancing the science and practice of conserving biodiversity. Conservation biology as a discipline reaches beyond biology, into subjects such as philosophy, law, economics, humanities, arts, anthropology, and education. Within biology, conservation genetics and evolution are immense fields unto themselves, but these disciplines are of prime importance to the practice and profession of conservation biology.

[...] there are advocates and there are sloppy or dishonest scientists, and these groups differ.

Chan

Is conservation biology an objective science when biologists advocate for an inherent value in nature? Do conservationists introduce bias when they support policies using qualitative description, such as habitat *degradation*, or *healthy* ecosystems? As all scientists hold values, so do conservation biologists. Conservation biologists advocate for reasoned and sensible management of natural resources and do so with a disclosed combination of science, reason, logic, and values in their conservation management plans. This sort of advocacy is similar to the medical profession advocating for healthy lifestyle options, both are beneficial to human well-being yet remain scientific in their approach. Many conservation biologists, in addition to having a Bachelors of Science (or extensive natural experience) often receive professional accreditation during their career.

There is a movement in conservation biology suggesting a new form of leadership is needed to mobilize conservation biology into a more effective discipline that is able to

communicate the full scope of the problem to society at large. The movement proposes an adaptive leadership approach that parallels an adaptive management approach. The concept is based on a new philosophy or leadership theory steering away from historical notions of power, authority, and dominance. Adaptive conservation leadership is reflective and more equitable as it applies to any member of society who can mobilize others toward meaningful change using communication techniques that are inspiring, purposeful, and collegial. Adaptive conservation leadership and mentoring programs are being implemented by conservation biologists through organizations such as the Aldo Leopold Leadership Program

## Approaches

Conservation may be classified as either in-situ conservation, which is protecting an endangered species in its natural habitat, or ex-situ conservation, which occurs outside the natural habitat. In-situ conservation involves protecting or cleaning up the habitat itself which may include a great deal of environmental preservation, or by defending the species from predators. Ex-situ conservation may be used on some or all of the population, when in-situ conservation is too difficult, or impossible.
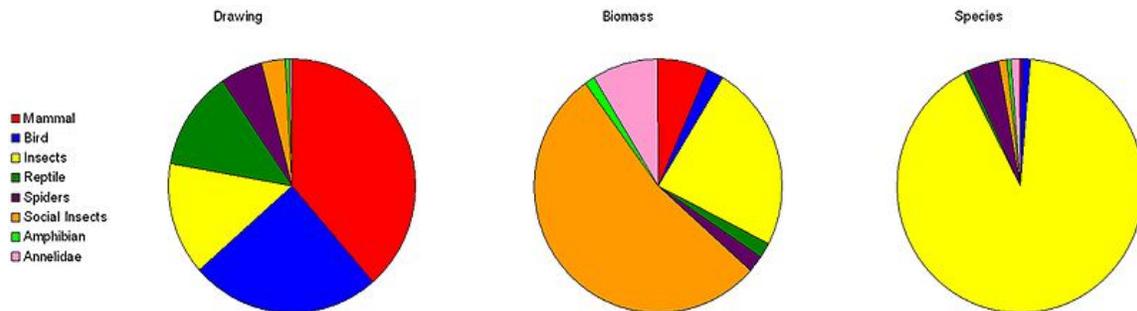
Also, non-interference may be used, which is termed a preservationist method. Preservationists advocate for giving areas of nature and species a protected existence that halts interference from the humans. In this regard, conservationists differ from preservationists in the social dimension, as conservation biology engages society and seeks equitable solutions for both society and ecosystems.

## Ethics and values

Conservation biologists are interdisciplinary researchers that practice ethics in the biological and social sciences. Chan states that conservationists must advocate for biodiversity and can do so in a scientifically ethical manner by not promoting simultaneous advocacy against other competing values. A conservationist researches biodiversity and reasons through a Resource Conservation Ethic, which identify what measures will deliver "the greatest good for the greatest number of people for the longest time.

Some conservation biologists argue that nature has an intrinsic value that is independent of anthropocentric usefulness or utilitarianism. Intrinsic value advocates that a gene, or species, be valued because they have a utility for the ecosystems they sustain. Aldo Leopold was a classical thinker and writer on such conservation ethics whose philosophy, ethics and writings are still valued and revisited by modern conservation biologists. His writing is oftentimes required reading for those in the profession.

## Conservation priorities



A pie chart image showing the relative biomass representation in a rain forest through a summary of children's perceptions from drawings and artwork (left), through a scientific estimate of actual biomass (middle), and by a measure of biodiversity (right). Notice that the biomass of social insects (middle) far outweighs the number of species (right).

While most in the community of conservation science "stress the importance" of sustaining biodiversity, there is debate on how to prioritize genes, species, or ecosystems, which are all components of biodiversity (e.g. Bowen, 1999). While the predominant approach to date has been to focus efforts on endangered species by conserving *biodiversity hotspots*, some scientists (e.g.) and conservation organizations, such as the Nature Conservancy, argue that it is more cost effective, logical, and socially relevant to invest in *biodiversity coldspots*. The costs of discovering, naming, and mapping out the distribution every species, they argue, is an ill advised conservation venture. They reason it is better to understand the significance of the ecological roles of species.

Biodiversity hotspots and coldspots are a way of recognizing that the spatial concentration of genes, species, and ecosystems is not uniformly distributed on the Earth's surface. For example, "[...] 44% of all species of vascular plants and 35% of all species in four vertebrate groups are confined to 25 hotspots comprising only 1.4% of the land surface of the Earth."

Those arguing in favor of setting priorities for coldspots point out that there are other measures to consider beyond biodiversity. They point out that emphasizing hotspots downplays the importance of the social and ecological connections to vast areas of the Earth's ecosystems where biomass, not biodiversity, reigns supreme. It is estimated that 36% of the Earth's surface, encompassing 38.9% of the worlds vertebrates, lacks the endemic species to qualify as biodiversity hotspot. Moreover, measures show that maximizing protections for biodiversity does not capture ecosystem services any better than targeting randomly chosen regions. Population level biodiversity (i.e. coldspots) are disappearing at a rate that is ten times that at the species level. The level of importance in addressing biomass versus endemism as a concern for conservation biology is highlighted in literature measuring the level of threat to global ecosystem carbon stocks that do not necessarily reside in areas of endemism. A hotspot priority approach would not invest so heavily in places such as steppes, the Serengeti, the Arctic, or taiga. These areas

contribute a great abundance of population (not species) level biodiversity and ecosystem services, including cultural value and planetary nutrient cycling.

Summary of 2006 IUCN Red List categories.

Those in favor of the hotspot approach point out that species are irreplaceable components of the global ecosystem, they are concentrated in places that are most threatened, and should therefore receive maximal strategic protections. This is a hotspot approach because the priority is set to target species level concerns over population level or biomass. Species richness and genetic biodiversity contributes to and engenders ecosystem stability, ecosystem processes, evolutionary adaptability, and biomass. Both sides agree, however, that conserving biodiversity is necessary to reduce the extinction rate and identify an inherent value in nature; the debate hinges on how to prioritize limited conservation resources in the most cost effective way.

## Economic values and natural capital



Tadrart Acacus desert in western Libya, part of the Sahara.

Conservation biologists have started to collaborate with leading global economists to determine how to measure the wealth and services of nature and to make these values apparent in global market transactions. This system of accounting is called *natural*

*capital* and would, for example, register the value of an ecosystem before it is cleared to make way for development. The WWF publishes its *Living Planet Report* and provides a global index of biodiversity by monitoring approximately 5,000 populations in 1,686 species of vertebrate (mammals, birds, fish, reptiles, and amphibians) and report on the trends in much the same way that the stock market is tracked.

This method of measuring the global economic benefit of nature has been endorsed by the G8+5 leaders and the European Commission. Nature sustains many ecosystem services that benefit humanity. Many of the earths ecosystem services are public goods without a market and therefore no price or value. When the *stock market* registers a financial crisis, traders on Wall Street are not in the business of trading stocks for much of the planet's living natural capital stored in ecosystems. There is no natural stock market with investment portfolios into sea horses, amphibians, insects, and other creatures that provide a sustainable supply of ecosystem services that are valuable to society. The ecological footprint of society has exceeded the bio-regenerative capacity limits of the planet's ecosystems by about 30 percent, which is the same percentage of vertebrate populations that have registered decline from 1970 through 2005.

The ecological credit crunch is a global challenge. The Living Planet Report 2008 tells us that more than three quarters of the world's people live in nations that are ecological debtors – their national consumption has outstripped their country's biocapacity. Thus, most of us are propping up our current lifestyles, and our economic growth, by drawing (and increasingly overdrawing) upon the ecological capital of other parts of the world.
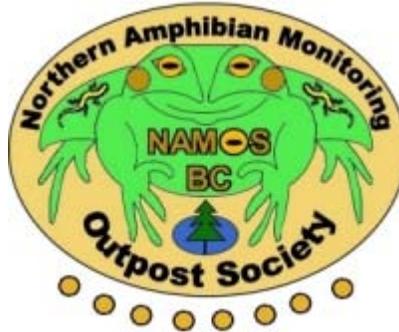
WWF Living Planet Report

The inherent natural economy plays an essential role in sustaining humanity, including the regulation of global atmospheric chemistry, pollinating crops, pest control, cycling soil nutrients, purifying our water supply, supplying medicines and health benefits, and unquantifiable quality of life improvements. There is a relationship, a correlation, between markets and natural capital, and social income inequity and biodiversity loss. This means that there are greater rates of biodiversity loss in places where the inequity of wealth is greatest

Although a direct market comparison of natural capital is likely insufficient in terms of human value, one measure of ecosystem services suggests the contribution amounts to trillions of dollars yearly. For example, one segment of North American forests has been assigned an annual value of 250 billion dollars; as another example, honey-bee pollination is estimated to provide between 10 and 18 billion dollars of value yearly. The value of ecosystem services on one New Zealand island has been imputed to be as great as the GDP of that region. This planetary wealth is being lost at an incredible rate as the demands of human society is exceeding the bio-regenerative capacity of the Earth. While biodiversity and ecosystems are resilient, the danger of losing them is that humans cannot recreate many ecosystem functions through technological innovation.

**Strategic species concepts**

 **Indicator species**



The NAMOS BC logo is an example of an ecosystem *umbrella* concept (forests and wetlands) combined with amphibians as *indicator* and *flagship species*.

An *indicator species* has a narrow set of ecological requirements, therefore they become useful targets for observing the health of an ecosystem. Some animals, such as amphibians with their semi-permeable skin and linkages to wetlands, have an acute sensitivity to environmental harm and thus may serve as a *miner's canary*. Indicator species are monitored in an effort to capture environmental degradation through pollution or some other link to proximate human activities. Monitoring an indicator species is a measure to determine if there is a significant environmental impact that can serve to advise or modify practice, such as through different forestsilviculture treatments and management scenarios, or to measure the degree of harm that a pesticide may impart on the health of an ecosystem.

Government regulators, consultants, or NGOs regularly monitor indicator species, however, there are limitations coupled with many practical considerations that must be followed for the approach to be effective. It is generally recommended that multiple indicators (genes, populations, species, communities, and landscape) be monitored for effective conservation measurement that prevents harm to the complex, and oftentimes unpredictable, response from ecosystem dynamics (Noss, 1997).

# Umbrella and flagship species

## Umbrella Spicies

**Umbrella species** are species selected for making conservation related decisions, typically because protecting these species indirectly protects the many other species that make up the ecological community of its habitat. Species conservation can be subjective because it is hard to determine the status of many species. With millions of species of concern, the identification of selected *keystone species*, *flagship species* or *umbrella*

*species* makes conservation decisions easier. Umbrella species can be used to help select the locations of potential reserves, find the minimum size of these conservation areas or reserves, and to determine the composition, structure and processes of ecosystems.

# Definitions

Two commonly used definitions:

- A: "A wide-ranging species whose requirements include those of many other species"
- B: A species with large area requirements for which protection of the species offers protection to other species that share the same habitat

Other descriptions include:

- A: "The protection of umbrella species automatically extends protection to other species. i.e. spotted owl and old growth trees"
- B: "Traditional umbrella species, relatively large-bodied and wide-ranging species of higher vertebrates"

# Use in landuse management

The use of umbrella species as a conservation tool is highly debated. The term was first used by Wilcox (1984) who defined an umbrella species as one whose minimum area requirements are at least as comprehensive of the rest of the community for which protection is sought though the establishment and management of a protected area.

Some scientists have found that the umbrella effect provides a simpler way to manage ecological communities. Others feel that a combination of other tools establish better land management reserves to help protect more species than just using umbrella species alone. Individual invertebrate species can be good umbrella species because they can protect older, unique ecosystems. There have been cases where umbrella species have protected a large amount of area which has been beneficial to surrounding species such as the northern spotted owl.

Currently research is being done on land management decisions based on using umbrella species to protect habitat of specific species as well as other organisms in the area. Dunk, Zielinski and Welsh (2006) reported that the reserves in Northern California (Klamath-Siskiyou forests), set aside for the northern spotted owl, also protect mollusks and salamanders within that habitat. According to their conclusions, the reserves set aside for the northern spotted owl "serve as a reasonable coarse-filter umbrella species for the taxa [they] evaluated," which were the mollusks and salamanders.

# Use in the Endangered Species Act (USA)

The Bay checkerspot butterfly has been on the Endangered Species List since 1987 and is still currently listed. Launer and Murphy (1994) tried to determine whether this butterfly could be considered an umbrella species in protecting the native grassland it inhabits. They discovered that the Endangered Species Act (ESA) has a loophole to eliminate federally protected plants that reside on private property. However, the California Environmental Quality Act (CEQA) reinforces state conservation regulations. Using the ESA to protect termed umbrella species and their habitats can be controversial because they are not as reinforced in some states as others (such as California) to protect overall biodiversity.

## Flagship Species

Project logo showing the use of the Zanzibar Red Colobus as the flagship species for conservation in Zanzibar

A **flagship species** is a species chosen to represent an environmental cause, such as an ecosystem in need of conservation. These species are chosen for their vulnerability,

attractiveness or distinctiveness in order to engender support and acknowledgment from the public at large. Thus, the concept of a flagship species holds that by giving publicity to a few key species, the support given to those species will successfully leverage conservation of entire ecosystems and all species contained therein.

Examples of flagship species include the Asiatic lion and the Bengal tiger of India, the giant panda of China, the golden lion tamarin of Brazil, the African elephant, the mountain gorilla of central Africa, and the orangutan of southeast Asia.

# History

**The conservation of natural resources is the fundamental problem. Unless we solve that problem, it will avail us little to solve all others.**

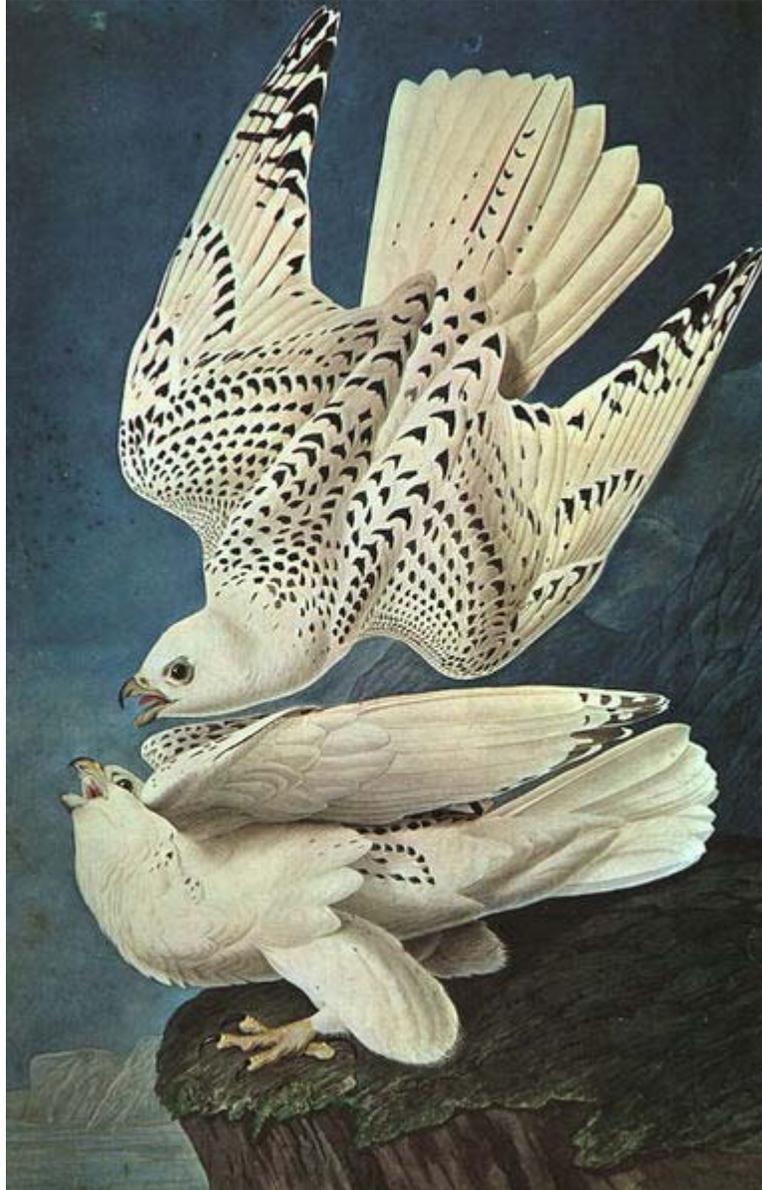Theodore Roosevelt

## Natural resource conservation

Efforts to conserve and protect *global* biodiversity are a recent phenomenon. Prior to the global conservation era, there was the coming of the age of conservation. Some historians have linked this with the 1916 National Parks Act, which included the 'use without impairment' clause, sought by John Muir. This eventually resulted in the removal of a proposal to build a dam in Dinosaur National Monument in 1959.

Natural resource conservation, however, has a history that extends prior to the age of conservation. Resource ethics grew out of necessity through direct relations with nature. Regulation or communal restraint became necessary to prevent selfish motives from taking more than could be locally sustained, therefore compromising the long-term supply for the rest of the community. This social dilemma with respect to natural resource management is often called the "Tragedy of the Commons". From this principal, conservation biologists can trace communal resource based ethics throughout cultures as a solution to communal resource conflict. For example, the Alaskan Tlingit peoples and the Haida of the Pacific Northwest had resource boundaries, rules, and restrictions among clans with respect to the fishing of Sockeye Salmon. These rules were guided by clan elders who knew life-long details of each river and stream they managed. There are numerous examples in history where cultures have followed rules, rituals, and organized practice with respect to communal natural resource management.

Conservation ethics are also found in early religious and philosophical writings. There are examples in the Tao, Shinto, Hindu, Islamic and Buddhist traditions. In Greek philosophy, Plato lamented about pasture land degradation: "What is left now is, so to say, the skeleton of a body wasted by disease; the rich, soft soil has been carried off and only the bare framework of the district left." In the bible, through Moses, God commanded to let the land rest from cultivation every seventh year. Before the 18th century, however, much of European culture considered it a pagan view to admire nature. Wilderness was denigrated while agricultural development was praised. However, as

early as AD 680 a wildlife sanctuary was founded on the Farne Islands by St Cuthbert in response to his religious beliefs.

**Early naturalists**

White Gerfalcons drawn by John James Audubon

Natural history was a major preoccupation in the 18th century, with grand expeditions and the opening of popular public displays in Europe and North America. By 1900 there were 150 natural history museums in Germany, 250 in Great Britain, 250 in the United States, and 300 in France. Preservationist or conservationist sentiments are a development in the late 18th to early 20th century. The 19th century fascination with natural history engendered a fervor to be the first to collect rare specimens with the goal

of doing so before they became extinct by other such collectors. Although his artistic works and romantic depiction of avian life inspired many bird enthusiasts and conservation organizations, the writings of John James Audubon, by modern standards, show insensitivity toward bird conservation as he shot and collected hundreds of specimens. Inspired by him, however, the first chapter of the Audubon Society started in 1905 for the purpose of protecting birds.

## Coming of the Age of Conservation

The modern concept of ecosystem services can be found in the late 19th century. "The utility of Natural History or its applicability to promote the material wealth of the State cannot be doubted. It was a great mistake to suppose that the subjects of Zoology, Botany, and Geology did not involve much that affects our comfort, convenience, health and wealth."

In the department of Woods and Forestry we should teach on the principals of conservation and teach on the lessons of economy rather than of waste in the natural resources of our country.

American Museum of Natural History, 1909

By the early 1800s biogeography was ignited through efforts of Alexander von Humboldt, DeCandolle, Lyell and Darwin; their efforts, while important in relating species to their environments, were part of the naturalist tradition and fell short of conservation biology proper. Darwin, for example, hunted and shot birds and kept natural history cabinets in line with Victorian tradition.

Modern roots of conservation biology can be found in the late 19th century Enlightenment period particularly in England and Scotland. A number of thinkers, among them notably Lord Monboddo, described the importance of "preserving nature"; much of this early emphasis had its origins in Christian theology.

## 20th century conservation

In the 20th century, actions in the United Kingdom, United States, and Canada emphasized the protection of habitat areas pursuant to visions of such people as John Muir, Theodore Roosevelt, and Aldo Leopold. While the Canadian nor the United Kingdom governments did not pioneer the creation of National Parks as the United States did in the late 19th century, there were many far-sighted civil servants who were dedicated to wildlife conservation and of notable mention. Some of these historical figures include Charles Gordon Hewitt and James Harkin.

The term *conservation* came into use in the late 19th century and referred to the management, mainly for economic reasons, of such natural resources as timber, fish, game, topsoil, pastureland, and minerals. In addition it referred to the preservation of forests (forestry), wildlife (wildlife refuge), parkland, wilderness, and watersheds. Western Europe was the source of much 19th century progress for conservation biology,

particularly the British Empire with the Sea Birds Preservation Act 1869. However, the United States made contributions to this field starting with thinking of Thoreau and taking form with the Forest Act of 1891, John Muir's founding of the Sierra Club in 1892, the founding of the New York Zoological Society in 1895 and establishment of a series of national forests and preserves by Theodore Roosevelt from 1901 to 1909.

Not until the mid 20th century did efforts arise to target individual species for conservation, notably efforts in big cat conservation in South America led by the New York Zoological Society. In the early 20th century the New York Zoological Society was instrumental in developing concepts of establishing preserves for particular species and conducting the necessary conservation studies to determine the suitability of locations that are most appropriate as conservation priorities; the work of Henry Fairfield Osborn Jr., Carl E. Akeley, Archie Carr and Archie Carr III is notable in this era. Akeley for example, having led expeditions to the Virunga Mountains and observed the mountain gorilla in the wild, became convinced that the species and the area were conservation priorities. He was instrumental in persuading Albert I of Belgium to act in defense of the mountain gorilla and establish Albert National Park (since renamed Virunga National Park) in what is now Democratic Republic of Congo.

By the 1970s, led primarily by work in the United States under the Endangered Species Act along with the Species at Risk Act (SARA) of Canada, Biodiversity Action Plans developed in Australia, Sweden, the United Kingdom, hundreds of species specific protection plans ensued. Notably the United Nations acted to conserve sites of outstanding cultural or natural importance to the common heritage of mankind. The programme was adopted by the General Conference of UNESCO in 1972. As of 2006, a total of 830 sites are listed: 644 cultural, 162 natural. The first country to pursue aggressive biological conservation through national legislation was the United States, which passed back to back legislation in the Endangered Species Act (1966) and National Environmental Policy Act (1970), which together injected major funding and protection measures to large scale habitat protection and threatened species research. Other conservation developments, however, have taken hold throughout the world. India, for example, passed the Wildlife Protection Act of 1972.

In 1980 a significant development was the emergence of the urban conservation movement. A local organization was established in Birmingham, UK, a development followed in rapid succession in cities across the UK, then overseas. Although perceived as a grassroots movement, its early development was driven by academic research into urban wildlife. Initially perceived as radical, the movement's view of conservation being inextricably linked with other human activity has now become mainstream in conservation thought. Considerable research effort is now directed at urban conservation biology. The Society for Conservation Biology originated in 1985.

By 1992 most of the countries of the world had become committed to the principles of conservation of biological diversity with the Convention on Biological Diversity; subsequently many countries began programmes of Biodiversity Action Plans to identify and conserve threatened species within their borders, as well as protect associated

habitats. The late 1990s saw increasing professionalism in the sector, with the maturing of organisations such as the Institute of Ecology and Environmental Management and the Society for the Environment.

Since 2000 the concept of landscape scale conservation has risen to prominence, with less emphasis being given to single-species or even single-habitat focused actions. Instead an ecosystem approach is advocated by most mainstream conservationist, although concerns have been expressed by those working to protect some high-profile species.
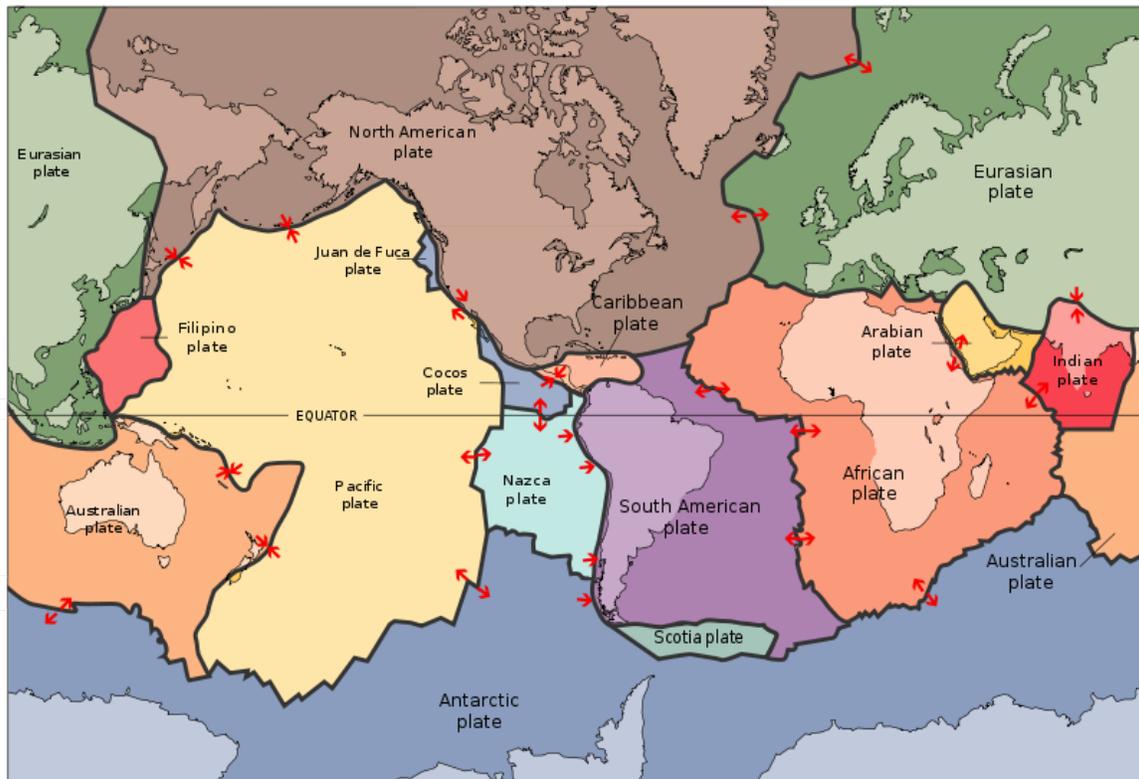
Ecology has clarified the workings of the biosphere; i.e., the complex interrelationships among humans, other species, and the physical environment. The burgeoning human population and associated agriculture, industry, and the ensuing pollution, have demonstrated how easily ecological relationships can be disrupted.

" The last word in ignorance is the man who says of an animal or plant: "What good is it?" If the land mechanism as a whole is good, then every part is good, whether we understand it or not. If the biota, in the course of aeons, has built something we like but do not understand, then who but a fool would discard seemingly useless parts? To keep every cog and wheel is the first precaution of " intelligent tinkering.

—Aldo Leopold, *A Sand County Almanac*

# Chapter-3

# Plate Tectonics



The tectonic plates of the world were mapped in the second half of the 20th century.

**Plate tectonics** (from the Late Latin *tectonicus*, from the Greek: τεκτονικός "pertaining to building") (Little, Fowler & Coulson 1990) is a scientific theory which describes the large scale motions of Earth's lithosphere. The theory builds on the older concepts of continental drift, developed during the first decades of the 20th century (one of the most famous advocates was Alfred Wegener), and was accepted by the majority of the Geoscientific community when the concepts of seafloor spreading were developed in the late 1950s and early 1960s. The lithosphere is broken up into what are called "tectonic plates". In the case of the Earth, there are currently seven to eight major (depending on how they are defined) and many minor plates. The lithospheric plates ride on the

asthenosphere. These plates move in relation to one another at one of three types of plate boundaries: convergent, or collisional boundaries; divergent boundaries, also called spreading centers; and conservative transform boundaries. Earthquakes, volcanic activity, mountain-building, and oceanic trench formation occur along these plate boundaries. The lateral relative movement of the plates varies, though it is typically 0–100 mm annually (Read & Watson 1975).

The tectonic plates are composed of two types of lithosphere: thicker continental and thin oceanic. The upper part is called the crust, again of two types (continental and oceanic). This means that a plate can be of one type, or of both types. One of the main points the theory proposes is that the amount of surface of the (continental and oceanic) plates that disappear in the mantle along the convergent boundaries by subduction is more or less in equilibrium with the new (oceanic) crust that is formed along the divergent margins by seafloor spreading. This is also referred to as the "conveyor belt" principle. In this way, the total surface of the Globe remains the same. This is in contrast with earlier theories advocated before the Plate Tectonics "paradigm", as it is sometimes called, became the main scientific model, theories that proposed gradual shrinking (contraction) or gradual expansion of the Globe, and that still exist in science as alternative models.

Regarding the driving mechanism of the plates various models co-exist: Tectonic plates are able to move because the Earth's lithosphere has a higher strength and lower density than the underlying asthenosphere. Lateral density variations in the mantle result in convection. Their movement is thought to be driven by a combination of the motion of seafloor away from the spreading ridge (due to variations in topography and density of the crust that result in differences in gravitational forces) and drag, downward suction, at the subduction zones. A different explanation lies in different forces generated by the rotation of the Globe and tidal forces of the Sun and the Moon. The relative importance of each of these factors is unclear.

# Key principles

The outer layers of the Earth are divided into lithosphere and asthenosphere. This is based on differences in mechanical properties and in the method for the transfer of heat. Mechanically, the lithosphere is cooler and more rigid, while the asthenosphere is hotter and flows more easily. In terms of heat transfer, the lithosphere loses heat by conduction whereas the asthenosphere also transfers heat by convection and has a nearly adiabatic temperature gradient. This division should not be confused with the *chemical* subdivision of these same layers into the mantle (comprising both the asthenosphere and the mantle portion of the lithosphere) and the crust: a given piece of mantle may be part of the lithosphere or the asthenosphere at different times, depending on its temperature and pressure.

The key principle of plate tectonics is that the lithosphere exists as separate and distinct *tectonic plates*, which ride on the fluid-like (visco-elastic solid) asthenosphere. Plate motions range up to a typical 10–40 mm/a (Mid-Atlantic Ridge; about as fast as fingernails grow), to about 160 mm/a (Nazca Plate; about as fast as hair grows) (Zhen
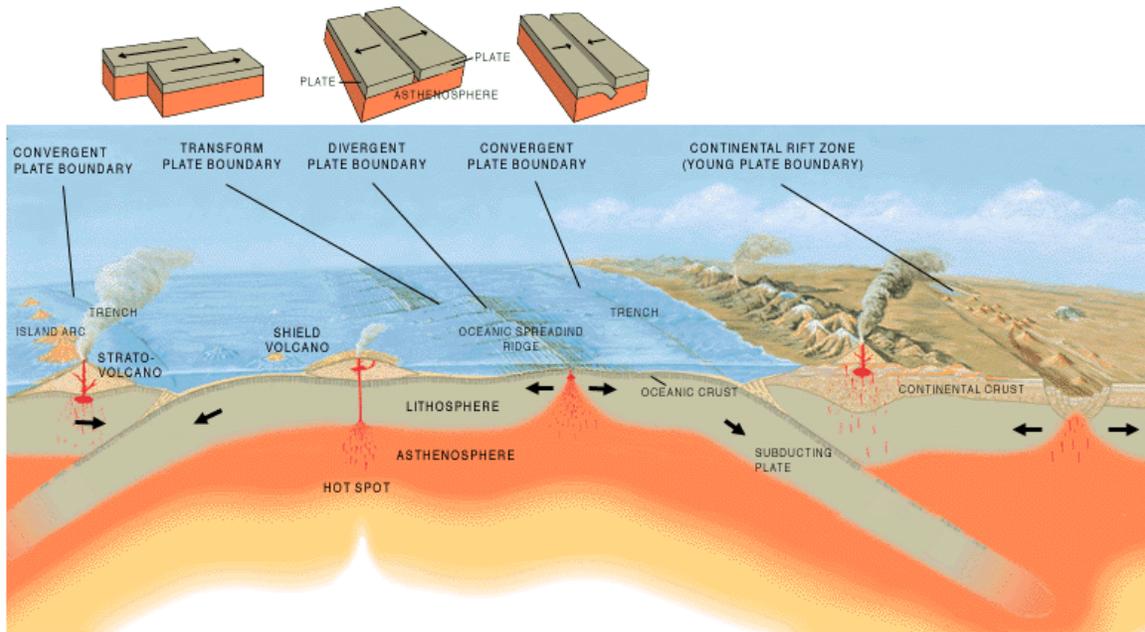
Shao 1997; Hancock, Skinner & Dineley 2000). The driving mechanism behind this movement is described below in a separate section.

Tectonic lithosphere plates consist of lithospheric mantle overlain by either or both of two types of crustal material: oceanic crust (in older texts called *sima* from silicon and magnesium) and continental crust (*sial* from silicon and aluminium). Average oceanic lithosphere is typically 100 km thick (Turcotte & Schubert 2002); its thickness is a function of its age: as time passes, it conductively cools and becomes thicker. Because it is formed at mid-ocean ridges and spreads outwards, its thickness is therefore a function of its distance from the mid-ocean ridge where it was formed. For a typical distance oceanic lithosphere must travel before being subducted, the thickness varies from about 6 km thick at mid-ocean ridges to greater than 100 km at subduction zones; for shorter or longer distances, the subduction zone (and therefore also the mean) thickness becomes smaller or larger, respectively (Turcotte & Schubert 2002). Continental lithosphere is typically ~200 km thick, though this also varies considerably between basins, mountain ranges, and stable cratonic interiors of continents. The two types of crust also differ in thickness, with continental crust being considerably thicker than oceanic (35 km vs. 6 km) (Turcotte & Schubert 2002).

The location where two plates meet is called a *plate boundary*, and plate boundaries are commonly associated with geological events such as earthquakes and the creation of topographic features such as mountains, volcanoes, mid-ocean ridges, and oceanic trenches. The majority of the world's active volcanoes occur along plate boundaries, with the Pacific Plate's Ring of Fire being most active and most widely known. These boundaries are discussed in further detail below.

As explained above, tectonic plates can include continental crust or oceanic crust, and many plates contain both. For example, the African Plate includes the continent and parts of the floor of the Atlantic and Indian Oceans. The distinction between oceanic crust and continental crust is based on their modes of formation. Oceanic crust is formed at sea-floor spreading centers, and continental crust is formed through arc volcanism and accretion of terranes through tectonic processes; though some of these terranes may contain ophiolite sequences, which are pieces of oceanic crust, these are considered part of the continent when they exit the standard cycle of formation and spreading centers and subduction beneath continents. Oceanic crust is also denser than continental crust owing to their different compositions. Oceanic crust is denser because it has less silicon and more heavier elements ("mafic") than continental crust ("felsic") (Schmidt & Harbert 1998). As a result of this density stratification, oceanic crust generally lies below sea level (for example most of the Pacific Plate), while the continental crust buoyantly projects above sea level.

# Types of plate boundaries



Three types of plate boundary.

Basically, three types of plate boundaries exist (Meissner 2002, p. 100), with a fourth, mixed type, characterized by the way the plates move relative to each other. They are associated with different types of surface phenomena. The different types of plate boundaries are:

1. *Transform boundaries (Conservative)* occur where plates slide or, perhaps more accurately, grind past each other along transform faults. The relative motion of the two plates is either sinistral (left side toward the observer) or dextral (right side toward the observer). The San Andreas Fault in California is an example of a transform boundary exhibiting dextral motion.
2. *Divergent boundaries (Constructive)* occur where two plates slide apart from each other. Mid-ocean ridges (e.g., Mid-Atlantic Ridge) and active zones of rifting (such as Africa's Great Rift Valley) are both examples of divergent boundaries.
3. *Convergent boundaries (Destructive)* (or *active margins*) occur where two plates slide towards each other commonly forming either a subduction zone (if one plate moves underneath the other) or a continental collision (if the two plates contain continental crust). Deep marine trenches are typically associated with subduction zones. The subducting slab contains many hydrous minerals, which release their water on heating; this water then causes the mantle to melt, producing volcanism. Examples of this are the Andes mountain range in South America and the Japanese island arc.
4. *Plate boundary zones* occur where the effects of the interactions are unclear and the broad belt boundaries are not well defined.

# Driving forces of plate motion

Plate tectonics is basically a kinematic phenomenon: Earth scientists agree upon the observation and deduction that the plates have moved one with respect to the other, and debate and find agreements on how and when. But still a major question remains on what is the motor behind this movement; the geodynamic mechanism, and here science diverges in different theories.

Generally, it is accepted that tectonic plates are able to move because of the relative density of oceanic lithosphere and the relative weakness of the asthenosphere. Dissipation of heat from the mantle is acknowledged to be the original source of energy driving plate tectonics, through convection or large scale upwelling and doming. As a consequence, in the current view, although it is still a matter of some debate, because of the excess density of the oceanic lithosphere sinking in subduction zones a powerful source of plate motion is generated. When the new crust forms at mid-ocean ridges, this oceanic lithosphere is initially less dense than the underlying asthenosphere, but it becomes denser with age, as it conductively cools and thickens. The greater density of old lithosphere relative to the underlying asthenosphere allows it to sink into the deep mantle at subduction zones, providing most of the driving force for plate motions. The weakness of the asthenosphere allows the tectonic plates to move easily towards a subduction zone. Although subduction is believed to be the strongest force driving plate motions, it cannot be the only force since there are plates such as the North American Plate which are moving, yet are nowhere being subducted. The same is true for the enormous Eurasian Plate. The sources of plate motion are a matter of intensive research and discussion among earth scientists. One of the main points is that the kinematic pattern of the movements itself should be separated clearly from the possible geodynamic mechanism that is invoked as the driving force of the observed movements, as some patterns may be explained by more than one mechanism (van Dijk 1992, van Dijk & Okkes 1991). Basically, the driving forces that are advocated at the moment, can be divided in three categories: mantle dynamics related, gravity related (mostly secondary forces), and Earth rotation related.

## Mantle dynamics related driving forces

For a considerable period of around 25 years (last quarter of the twentieth century) the leading theory envisaged large scale convection currents in the upper mantle which are transmitted through the asthenosphere as the main driving force of the tectonic plates. This theory was launched by Arthur Holmes and some forerunners in the 1930s and was immediately recognised as the solution for the acceptance of the theory discussed since its occurrence in the papers of Alfred Wegener in the early years of the century. It was, though, long debated because the leading ("fixist") theory was still envisaging a static Earth without moving continents, up until the major break–throughs in the early sixties.

Two– and three–dimensional imaging of the Earth's interior (seismic tomography) shows that there is a laterally varying density distribution throughout the mantle. Such density variations can be material (from rock chemistry), mineral (from variations in mineral structures), or thermal (through thermal expansion and contraction from heat energy).

The manifestation of this varying lateral density is mantle convection from buoyancy forces (Tanimoto & Lay 2000).

How mantle convection relates directly and indirectly to the motion of the plates is a matter of ongoing study and discussion in geodynamics. Somehow, this energy must be transferred to the lithosphere in order for tectonic plates to move. There are essentially two types of forces that are thought to influence plate motion: friction and gravity.

Basal drag (friction): The plate motion is in this way driven by friction between the convection currents in the asthenosphere and the more rigid overlying floating lithosphere.
Slab suction (gravity): Local convection currents exert a downward frictional pull on plates in subduction zones at ocean trenches. Slab suction may occur in a geodynamic setting wherein basal tractions continue to act on the plate as it dives into the mantle (although perhaps to a greater extent acting on both the under and upper side of the slab).

Lately, the convection theory is much debated as modern techniques based on 3D seismic tomography of imaging the internal structure of the Earth's mantle still fail to recognise these predicted large scale convection cells. Therefore, alternative patterns of mantle dynamics have been proposed:

In the theory of plume tectonics developed during the 1990s, a modified concept of mantle convection currents is used, related to super plumes rising from the deeper mantle which would be the drivers or the substitutes of the major convection cells. These ideas, which find their roots in the early 1930s with the so-called "fixistic" ideas of the European and Russian Earth Science Schools, find resonance in the modern theories which envisage hot spots/mantle plumes in the mantle which remain fixed and are overridden by oceanic and continental lithosphere plates during time, and leave their traces in the geological record (though these phenomena are not invoked as real driving mechanisms, but rather as a modulator). The modern theories that continue building on the older mantle doming concepts and see the movements of the plates a secondary phenomena, are beyond the scope of this page and are discussed elsewhere for example on the plume tectonics page.
Another suggestion is that the mantle flows neither in cells nor large plumes, but rather as a series of channels just below the Earth's crust which then provide basal friction to the lithosphere. This theory is called "surge tectonics" and became quite popular in geophysics and geodynamics during the 1980s and 1990s (Smoot et al. 1996).

## Gravity related driving forces

Gravity related forces are usually invoked as secondary phenomena within the framework of a more general driving mechanism such as the various forms of mantle dynamics described above.

Gravitational sliding away from a spreading ridge: According to many authors, plate motion is driven by the higher elevation of plates at ocean ridges. As oceanic lithosphere

is formed at spreading ridges from hot mantle material, it gradually cools and thickens with age (and thus distance from the ridge). Cool oceanic lithosphere is significantly denser than the hot mantle material from which it is derived and so with increasing thickness it gradually subsides into the mantle to compensate the greater load. The result is a slight lateral incline with distance from the ridge axis.

This force is regarded as a secondary force often referred to as "ridge-push". This is a misnomer as nothing is "pushing" horizontally and tensional features are dominant along ridges. It is more accurate to refer to this mechanism as gravitational sliding as variable topography across the totality of the plate can vary considerably and the topography of spreading ridges is only the most prominent feature. Other mechanisms generating this gravitational secondary force are for example:

Flexural bulging of the lithosphere before it dives underneath an adjacent plate, for instance, produces a clear topographical feature that can offset or at least affect the influence of topographical ocean ridges.
Mantle plumes and hot spots impinging on the underside of tectonic plates can drastically alter the topography of the ocean floor. Some of these, on a larger scale, are seen as the major driving force of the plates (see below).

Slab-pull: Current scientific opinion is that the asthenosphere is insufficiently competent or rigid to directly cause motion by friction along the base of the lithosphere. Slab pull is therefore most widely thought to be the greatest force acting on the plates. In this current understanding, plate motion is mostly driven by the weight of cold, dense plates sinking into the mantle at trenches (Conrad & Lithgow-Bertelloni 2002). Recent models indicate that trench suction plays an important role as well. However, as the North American Plate is nowhere being subducted, yet it is in motion presents a problem. The same holds for the African, Eurasian, and Antarctic plates. Slab pull is especially invoked in areas where remnants of older lithosphere become trapped along convergence zones e.g. as relicts in collisional belts, which, sinking into the mantle and rolling backwards, exert a pull on the overlying crust.

Gravitational sliding away from mantle doming: According to older theories one of the driving mechanisms of the plates is the existence of large scale asthenosphere/mantle domes, which cause the gravitational sliding of lithosphere plates away from them. This gravitational sliding represents a secondary phenomenon of this, basically vertically oriented mechanism. This can act on various scales, from the small scale of one island arc up to the larger scale of an entire ocean basin.

## Earth rotation related driving forces

Alfred Wegener, being a meteorologist, had proposed tidal forces and pole flight Force as main driving mechanisms for continental drift. However, these forces were considered far too small to cause continental motion as the concept then was of continents plowing through oceanic crust. Therefore, also Wegener in his last edition of his book in 1929 converted to convection currents as the main driving force.

In the plate tectonics context (accepted since the seafloor spreading proposals of Heezen, Hess, Dietz, Morley, Vine and Matthews -see below- during the early 1960s), though, oceanic crust in motion *with* the continents which made the proposals related to Earth rotation to be reconsidered, also in more recent literature, these are:

1. Tidal drag due to the gravitational force the Moon (and the Sun) exerts on the crust of the Earth
2. Shear strain of the Earth globe due to N-S compression related to the rotation and modulations of it
3. Pole flight force: equatorial drift due to rotation and centrifugal effects: tendency of the plates to move from the poles to the equator ("*Polflucht*")
4. Coriolis effect the plates suffer when they move around the globe (coriolis effect/law of Buys Ballot)
5. Global deformation of the geoid due to small displacements of rotational pole with respect to the Earth crust
6. Other smaller deformation effects of the crust due to wobbles and spin movements of the Earth rotation on a smaller time scale.

In order for these mechanisms to be overall valid, systematic relationships should exist all over the Globe between the orientation and kinematics of deformation, and the geographical latitudinal and longitudinal grid of the Earth itself. Ironically, these systematic relations studies in the second half of the nineteenth century and the first half of the twentieth century do underline exactly the opposite: that the plates had not moved in time, that the deformation grid was fixed with respect to the Earth equator and axis, and that gravitational driving forces were generally acting vertically and caused only locally horizontal movements (the so-called pre-plate tectonic, "fixist theories"). Later studies (discussed below on this page) therefore invoked many of the relationships recognised during this pre-plate tectonics period, to support their theories.

Of the many forces discussed in this paragraph, tidal force is still highly debated and defended as a possible principle driving force, whereas the other forces are used or in global geodynamic models not using the plate tectonics concepts (therefore beyond the discussions treated in this section), or proposed as minor modulations within the overall plate tectonics model.

In 1973 George W. Moore of the USGS and R. C. Bostrom presented evidence for a general westward drift of the Earth's lithosphere with respect to the mantle, and, therefore, tidal forces or tidal lag or "friction" due to the Earth's rotation and the forces acting upon it by the Moon being a driving force for plate tectonics: as the Earth spins eastward beneath the moon, the moon's gravity ever so slightly pulls the Earth's surface layer back westward, just like proposed by Alfred Wegener (see above). In a more recent 2006 study (Scoppola et al. 2006), scientists rediscussed and advocated these earlier proposed ideas. It has also been suggested recently in Lovett (2006) that this observation may also explain why Venus and Mars have no plate tectonics, since Venus has no moon and Mars' moons are too small to have significant tidal effects on Mars. In a recent paper by Torsvik et al. (2010) it was suggested that, on the other hand, it can easily be observed

that many plates are moving north and eastward, and that the dominantly westward motion of the Pacific ocean basins derives simply from the eastward bias of the Pacific spreading center (which is not a predicted manifestation of such lunar forces). In the same paper the authors admit, however, that relative to the lower mantle, there is a slight westward component in the motions of all the plates. They demonstrated though that the westward drift, seen only for the past 30 Ma, is attributed to the increased dominance of the steadily growing and accelerating Pacific plate. The debate is still open.

**Relative significance of each driving force mechanism**

The actual vector of a plate's motion must necessarily be a function of all the forces acting upon the plate. However, therein remains the problem regarding what degree each process contributes to the motion of each tectonic plate.
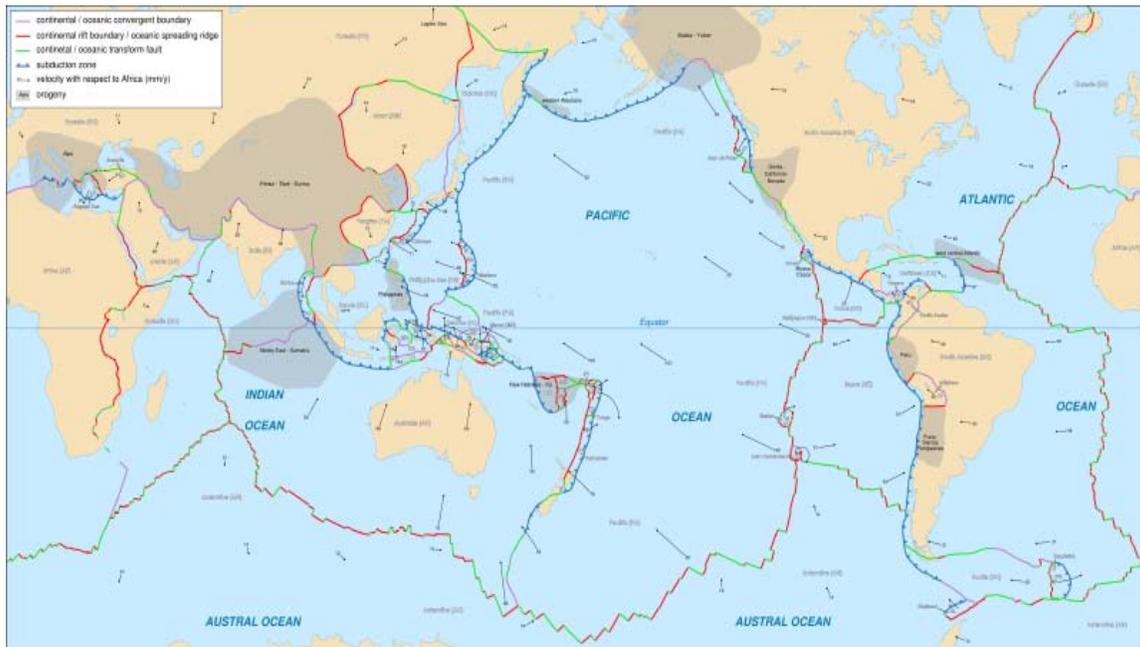
The diversity of geodynamic settings and properties of each plate must clearly result in differences in the degree to which such processes are actively driving the plates. One method of dealing with this problem is to consider the relative rate at which each plate is moving and to consider the available evidence of each driving force upon the plate as far as possible.

One of the most significant correlations found is that lithospheric plates attached to downgoing (subducting) plates move much faster than plates not attached to subducting plates. The Pacific plate, for instance, is essentially surrounded by zones of subduction (the so-called Ring of Fire) and moves much faster than the plates of the Atlantic basin, which are attached (perhaps one could say 'welded') to adjacent continents instead of subducting plates. It is thus thought that forces associated with the downgoing plate (slab pull and slab suction) are the driving forces which determine the motion of plates, except for those plates which are not being subducted (Conrad & Lithgow-Bertelloni 2002). The driving forces of plate motion continue to be active subjects of on-going research within geophysics and tectonophysics.

# Historical context - development of the theory

Plate tectonics is the main current theory in Earth Sciences regarding the development of our planet Earth. It is, therefore, appropriate to dedicate some space to explain how the Earth Science community, step by step, has built this theory, from early speculations, through the gathering of proof and severe debates, up to the refinement and quantification, and still ongoing confrontations with alternative ideas.

# Summary

Detailed map showing the tectonic plates with their movement vectors.

In line with other previous and contemporaneous proposals, in 1912 the meteorologist Alfred Wegener amply described what he called continental drift, expanded in his 1915 book *The Origin of Continents and Oceans* and the scientific debate started that would end up fifty years later in the theory of plate tectonics (Hughes 2001a). Starting from the idea (also expressed by his forerunners) that the present continents once formed a single land mass (which was called Pangea later on) that drifted apart, thus releasing the continents from the Earth's mantle and likening them to "icebergs" of low density granite floating on a sea of denser basalt (Wegener 1966; Hughes 2001b).

But without detailed evidence and a force sufficient to drive the movement, the theory was not generally accepted: the Earth might have a solid crust and mantle and a liquid core, but there seemed to be no way that portions of the crust could move around.

Notwithstanding much opposition, the view of continental drift gained support and a lively debate started between "drifters" or "mobilists" (proponents of the theory) and "fixists" (opponents). During the 1920s, 1930s and 1940s, the former reached important milestones proposing that convection currents might have driven the plate movements, and that spreading may have occurred below the sea within the oceanic crust. Concepts close to the elements now incorporated in plate tectonics were proposed by geophysisists and geologists (both fixists and mobilists) like Vening-Meinesz, Holmes, and Umbgrove.

One of the first pieces of geophysical evidence that was used to support the movement of lithospheric plates came from paleomagnetism. This is based on the fact that rocks of

different ages show a variable magnetic field direction, evidenced by studies since the mid–nineteenth century. The magnetic north and south reverse through time, and, especially important in paleotectonic studies, the relative position of the magnetic north varies through time. Initially, during the first half of the twentieth century, the latter phenomena was explained by introducing what was called "polar wander", i.e., it was assumed that the north pole location had been shifting through time. An alternative explanation, though, was that the continents had moved (shifted and rotated) relative to the north pole, and each continent, in fact, shows its own "polar wander path". During the late 1950s in was shown with success that these data could show the validity of continental drift in two occasions: by Keith Runcorn in a paper in 1956, and by Warren Carey in a symposium held in March 1956.

The second piece of evidence in support of continental drift came during the late 1950s and early 60s from data on the bathymetry of the deep ocean floors and the nature of the oceanic crust such as magnetic properties and, more generally, with the development of marine geology which gave evidence for the association of seafloor spreading along the mid-oceanic ridges and magnetic field reversals, published between 1959 and 1963 by Heezen, Dietz, Hess, Mason, Vine & Matthews, and Morley (Korgen 1995; Spiess & Kuperman 2003).

Simultaneous advances in early seismic imaging techniques in and around Wadati-Benioff zones along the trenches bounding many continental margins, together with many other geophysical (e.g. gravimetric) and geological observations, showed how the oceanic crust could disappear into the mantle, providing the mechanism to balance the extension of the ocean basins with shortening along its margins.

All these evidences, both from the ocean floor and from the continental margins made clear around 1965 that continental drift was feasible and the theory of plate tectonics, which was defined in a series of papers between 1965 and 1967, was born, with all its extraordinary explanatory and predictive power. The theory revolutionized the Earth sciences, explaining a diverse range of geological phenomena and their implications in other studies such as paleogeography and paleobiology.

## Continental drift

In the late 19th and early 20th centuries, geologists assumed that the Earth's major features were fixed, and that most geologic features such as basin development and mountain ranges could be explained by vertical crustal movement, described in what is called the geosynclinal theory. Generally, this was placed in the context of a contracting planet Earth due to heat loss in the course of a relatively short geological time.

It was observed as early as 1596 that the opposite coasts of the Atlantic Ocean—or, more precisely, the edges of the continental shelves—have similar shapes and seem to have once fitted together (Kious & Tilling 1996).

Since that time many theories were proposed to explain this apparent complementarity, but the assumption of a solid Earth made these various proposals difficult to accept (Frankel 1987).

The discovery of radioactivity and its associated heating properties in 1895 prompted a re-examination of the apparent age of the Earth (Joly 1909) since this had previously been estimated by its cooling rate and assumption the Earth's surface radiated like a black body (Thomson 1863).

Those calculations had implied that, even if it started at red heat, the Earth would have dropped to its present temperature in a few tens of millions of years. Armed with the knowledge of a new heat source, scientists realized that the Earth would be much older, and that its core was still sufficiently hot to be liquid.

By 1915, after having published a first article in 1912 (Wegener 1912). Alfred Wegener was making serious arguments for the idea of continental drift in the first edition of *The Origin of Continents and Oceans*. In that book (re-issued in four successive editions up to the final one in 1936), he noted how the east coast of South America and the west coast of Africa looked as if they were once attached. Wegener wasn't the first to note this (Abraham Ortelius, Snider-Pellegrini, Roberto Mantovani and Frank Bursley Taylor preceded him just to mention a few), but he was the first to marshal significant fossil and paleo-topographical and climatological evidence to support this simple observation (and was supported in this by researchers such as Alex du Toit). Furthermore, when the rock strata of the margins of separate continents are very similar it suggests that these rocks were formed in the same way, implying that they were joined initially. For instance, some parts of Scotland and Ireland contain rocks very similar to those found in Newfoundland and New Brunswick. Furthermore, the Caledonian Mountains of Europe and parts of the Appalachian Mountains of North America are very similar in structure and lithology.

However, his ideas were not taken seriously by many geologists, who pointed out that there was no apparent mechanism for continental drift. Specifically, they did not see how continental rock could plow through the much denser rock that makes up oceanic crust. Wegener could not explain the force that drove continental drift, and his vindication did not come until after his death in 1930.

## Floating continents - paleomagnetism - seismicity zones

As it was observed early that although granite existed on continents, seafloor seemed to be composed of denser basalt, the prevailing concept during the first half of the twentieth century was that there were two types of crust, named "sial" (continental type crust), and "sima" (oceanic type crust). Furthermore, it was supposed that a static shells of strata was present under the continents. It therefore looked apparent that a layer of basalt (sial) underlies the continental rocks.

However, based upon abnormalities in plumb line deflection by the Andes in Peru, Pierre Bouguer had deduced that less-dense mountains must have a downward projection into

the denser layer underneath. The concept that mountains had "roots" was confirmed by George B. Airy a hundred years later during study of Himalayan gravitation, and seismic studies detected corresponding density variations. Therefore, by the mid–1950s the question remained unresolved of whether mountain roots were clenched in surrounding basalt or were floating upon it like an iceberg.

During the 20th century, improvements in and greater use of seismic instruments such as seismographs enabled scientists to learn that earthquakes tend to be concentrated in specific areas, most notably along the oceanic trenches and spreading ridges. By the late 1920s, seismologists were beginning to identify several prominent earthquake zones parallel to the trenches that typically were inclined 40–60° from the horizontal and extended several hundred kilometers into the Earth. These zones later became known as Wadati-Benioff zones, or simply Benioff zones, in honor of the seismologists who first recognized them, Kiyoo Wadati of Japan and Hugo Benioff of the United States. The study of global seismicity greatly advanced in the 1960s with the establishment of the Worldwide Standardized Seismograph Network (WWSSN) to monitor the compliance of the 1963 treaty banning above-ground testing of nuclear weapons. The much improved data from the WWSSN instruments allowed seismologists to map precisely the zones of earthquake concentration world wide.

Meanwhile, debates developed around the phenomena of polar Wander. Since the early debates of continental drift, scientists had discussed and used evidence that polar drift had occurred due to the fact that continents seemed to have moved through different climatic zones during the past. Furthermore, paleomagnetic data had shown that the magnetic pole had also shifted during time. Reasoning in an opposite way, the continents might have shifted and rotated, while the pole remained relatively fixed. The first time the evidence of magnetic polar wander was used to support the movements of continents was in a paper by Keith Runcorn in 1956, and successive papers by him and his students Ted Irving (who was actually the first to be convinced of the fact that paleomagnetism supported continental drift) and Ken Creer.

This was immediately followed by a symposium in Tasmania in March 1956 (Carey 1956; the evidence was used in the theory of an expansion of the global crust. In this hypothesis the shifting of the continents can be simply explained by a large increase in size of the Earth since its formation. However, this was unsatisfactory because its supporters could offer no convincing mechanism to produce a significant expansion of the Earth. Certainly there is no evidence that the moon has expanded in the past 3 billion years; other work would soon show that the evidence was equally in support of continental drift on a globe with a stable radius.

During the thirties up to the late fifties, numerous milestones were reached that would eventually lead to the development of plate tectonics. These are the works of Vening-Meinesz, Holmes, Umbgrove, and numerous others, in which concepts close or near identical to modern plate tectonics theory where defined and outlined. The most important milestone was reached when the English geologist Arthur Holmes proposed in

1920 that plate junctions might lie beneath the sea, and in 1928 that convection currents within the mantle might be the driving force.

Often, all these milestones are forgotten for various reasons:

1. During this timespan, continental drift was not accepted.
2. Some of these ideas were discussed in the context of abandoned fixistic ideas of a deforming globe without continental drift or an expanding Earth.
3. They were published during an episode of extreme political and economic instability and scientific communication was obviously hampered by this.
4. Many of these were published by European scientists and at first not mentioned or given little credit in the papers published by the American researchers which during the 1960s presented evidence for sea floor spreading.

## Mid oceanic ridge spreading and convection

In 1947, a team of scientists led by Maurice Ewing utilizing the Woods Hole Oceanographic Institution's research vessel *Atlantis* and an array of instruments, confirmed the existence of a rise in the central Atlantic Ocean, and found that the floor of the seabed beneath the layer of sediments consisted of basalt, not the granite which is the main constituent of continents. They also found that the oceanic crust was much thinner than continental crust. All these new findings raised important and intriguing questions (Lippsett 2001; and Lippsett 2006).

The new data that had been collected on the ocean basins also showed particular characteristics regarding the bathymetry. One of the major outcomes of these datasets was that all along the globe, a system of mid-oceanic ridges was detected. An important conclusion was that along this system, new ocean floor was being created, which led to the concept of the "Great Global Rift". This was described in the crucial paper of Bruce Heezen (1960) which would trigger a real revolution in thinking. A profound consequence of seafloor spreading is that new crust was, and is now, being continually created along the oceanic ridges. Therefore, Heezen advocated the so-called "expanding Earth" hypothesis of S. Warren Carey (see above). So, still the question remained: how can new crust be continuously added along the oceanic ridges without increasing the size of the Earth? In reality, this question had been solved already by numerous scientists during the forties and the fifties, like Arthur Holmes, Vening-Meinesz, Coates and many others: The crust in excess disappeared along what were called the oceanic trenches where so-called "subduction" occurred. Therefore, when various scientists during the early sixties started to reason on the data at their disposal regarding ocean floor, the pieces of the theory fell quickly at its place.
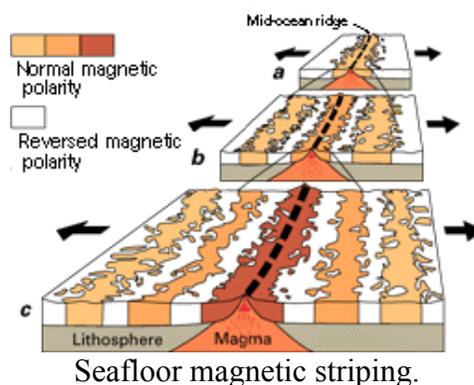
The question particularly intrigued Harry Hammond Hess, a Princeton University geologist and a Naval Reserve Rear Admiral, and Robert S. Dietz, a scientist with the U.S. Coast and Geodetic Survey who first coined the term *seafloor spreading*. Dietz and Hess (the former published the same idea one year earlier in *Nature*, but priority belongs to Hess who had already distributed an unpublished manuscript of his 1962 article by

1960) were among the small handful who really understood the broad implications of sea floor spreading and how it would eventually agree with the, at that time, unconventional and unaccepted ideas of continental drift and the elegant and mobilistic models proposed by previous workers like Holmes.

In the same year, Robert R. Coats of the U.S. Geological Survey described the main features of island arc subduction in the Aleutian Islands. His paper, though little–noted (and even ridiculed) at the time, has since been called "seminal" and "prescient". In reality, it actually shows that the work by the European scientists on island arcs and mountain belts performed and published during the 1930s up until the 1950s was applied and appreciated also in the United States.

If the Earth's crust was expanding along the oceanic ridges, Hess and Dietz reasoned like Holmes and others before them, it must be shrinking elsewhere. Hess followed Heezen suggesting that new oceanic crust continuously spreads away from the ridges in a conveyor belt–like motion. And, using the mobilistic concepts developed before, he correctly concluded that many millions of years later, the oceanic crust eventually descends along the continental margins where oceanic trenches – very deep, narrow canyons are present e.g. along the rim of the Pacific Ocean basin – were formed. The important step Hess made was that convection currents would be the driving force in this process, arriving at the same conclusions as Holmes had decades before with the only difference that the thinning of the ocean crust was performed using the mechanism of Heezen of spreading along the ridges. Hess therefore concluded that the Atlantic Ocean was expanding while the Pacific Ocean was shrinking. As old oceanic crust is "consumed" in the trenches, (like Holmes and others, he believed this was done by thickening of the continental lithosphere, not, as nowadays believed, by underthrusting at a larger scale of the oceanic crust itself into the mantle) new magma rises and erupts along the spreading ridges to form new crust. In effect, the ocean basins are perpetually being "recycled," with the creation of new crust and the destruction of old oceanic lithosphere occurring simultaneously, in a way like what later would be called the Wilson cycle (see below). Thus, the new mobilistic concepts neatly explained why the Earth does not get bigger with sea floor spreading, why there is so little sediment accumulation on the ocean floor, and why oceanic rocks are much younger than continental rocks.

### The final proof: magnetic striping



Seafloor magnetic striping.

A demonstration of magnetic striping. (The darker the color is the closer it is to normal polarity)

Beginning in the 1950s, scientists like Victor Vacquier, using magnetic instruments (magnetometers) adapted from airborne devices developed during World War II to detect submarines, began recognizing odd magnetic variations across the ocean floor. This finding, though unexpected, was not entirely surprising because it was known that basalt—the iron-rich, volcanic rock making up the ocean floor—contains a strongly magnetic mineral (magnetite) and can locally distort compass readings. This distortion was recognized by Icelandic mariners as early as the late 18th century. More important, because the presence of magnetite gives the basalt measurable magnetic properties, these newly discovered magnetic variations provided another means to study the deep ocean floor. When newly formed rock cools, such magnetic materials recorded the Earth's magnetic field at the time.

As more and more of the seafloor was mapped during the 1950s, the magnetic variations turned out not to be random or isolated occurrences, but instead revealed recognizable patterns. When these magnetic patterns were mapped over a wide region, the ocean floor showed a zebra-like pattern: one stripe with normal polarity and the adjoining stripe with reversed polarity. The overall pattern, defined by these alternating bands of normally and reversely polarized rock, became known as magnetic striping, and was published by Ron G. Mason and co-workers in 1961, who didn't find, though, an explanation for these data in terms of sea floor spreading, like Vine, Matthews and Morley a few years later (Mason & Raff 1961); (Raff & Mason 1961).

The discovery of magnetic striping called for an explanation. In the early 1960s scientists such as Heezen, Hess and Dietz had begun to theorise that mid-ocean ridges mark structurally weak zones where the ocean floor was being ripped in two lengthwise along the ridge crest. New magma from deep within the Earth rises easily through these weak zones and eventually erupts along the crest of the ridges to create new oceanic crust. This process, at first denominated the "conveyer belt hypothesis" and later called seafloor spreading, operating over many millions of years continues to form new ocean floor all across the 50,000 km-long system of mid–ocean ridges.

Only four years after the maps with the "zebra pattern" of magnetic stripes were published, the link between sea floor spreading and these patterns was correctly placed, independently by Lawrence Morley, and by Fred Vine and Drummond Matthews, in 1963 (Vine & Matthews 1963) now called the Vine-Matthews-Morley hypothesis. This

hypothesis linked these patterns to geomagnetic reversals and was supported by several lines of evidence:

1. the stripes are symmetrical around the crests of the mid-ocean ridges; at or near the crest of the ridge, the rocks are very young, and they become progressively older away from the ridge crest;
2. the youngest rocks at the ridge crest always have present-day (normal) polarity;
3. stripes of rock parallel to the ridge crest alternate in magnetic polarity (normal-reversed-normal, etc.), suggesting that they were formed during different epochs documenting the (already known from independent studies) normal and reversal episodes of the Earth's magnetic field.

By explaining both the zebra-like magnetic striping and the construction of the mid-ocean ridge system, the seafloor spreading hypothesis (SFS) quickly gained converts and represented another major advance in the development of the plate-tectonics theory. Furthermore, the oceanic crust now came to be appreciated as a natural "tape recording" of the history of the geomagnetic field reversals (GMFR) of the Earth's magnetic field. Nowadays, extensive studies are dedicated to the calibration of the normal-reversal patterns in the oceanic crust on one hand and known timescales derived from the dating of basalt layers in sedimentary sequences (magnetostratigraphy) on the other, to arrive at estimates of past spreading rates and plate reconstructions.

## Definition and refining of the theory - from new global tectonics to plate tectonics

After all these considerations, Plate Tectonics (or, as it was initially called "New Global Tectonics") became quickly accepted in the scientific world, and numerous papers followed that defined the concepts:
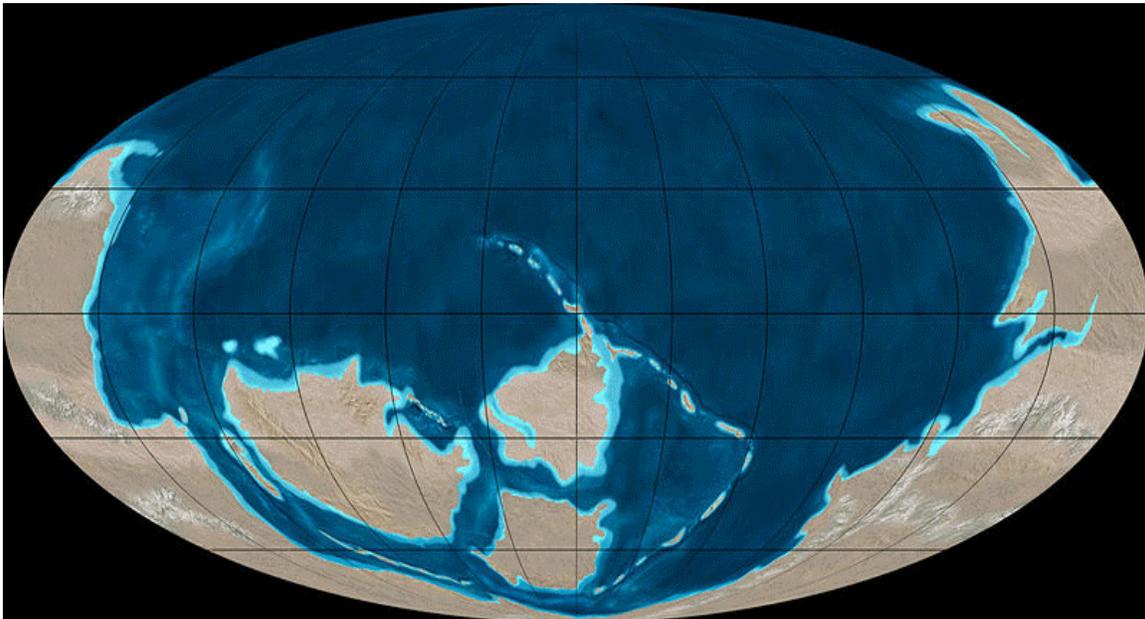
- In 1965, Tuzo Wilson who had been a promotor of the sea floor spreading hypothesis and continental drift from the very beginning (e.g. Wilson 1963) added the concept of transform faults to the model, completing the classes of fault types necessary to make the mobility of the plates on the globe work out (Wilson 1965).
- A symposium on continental drift was held at the Royal Society of London in 1965 which must be regarded as the official start of the acceptance of plate tectonics by the scientific community, and which abstracts are issued as Blacket, Bullard & Runcorn (1965). In this symposium, Edward Bullard and co-workers showed with a computer calculation how the continents along both sides of the Atlantic would best fit to close the ocean, which became known as the famous "Bullard's Fit".
- In 1966 Tuzo Wilson published the paper that referred to previous plate tectonic reconstructions, introducing the concept of what is now known as the "Wilson Cycle" (Wilson 1966).
- In 1967, at the American Geophysical Union's meeting, W. Jason Morgan proposed that the Earth's surface consists of 12 rigid plates that move relative to each other (Morgan 1968).

- Two months later, Xavier Le Pichon published a complete model based on 6 major plates with their relative motions, and we may say that this parks the final acceptance of the scientific community of plate tectonics (Le Pichon 1967).
- In the same year, McKenzie and Parker independently presented a model similar to Morgan's using translations and rotations on a sphere to define the plate motions (McKenzie & Parker 1967).

# Implications for biogeography

Continental drift theory helps biogeographers to explain the disjunct biogeographic distribution of present day life found on different continents but having similar ancestors (Moss & Wilson 1998). In particular, it explains the Gondwanan distribution of ratites and the Antarctic flora.

# Plate reconstruction



Reconstruction of plate configurations for the whole Phanerozoic

Reconstruction is used to establish past (and future) plate configurations, helping determine the shape and make-up of ancient supercontinents and providing a basis for paleogeography.
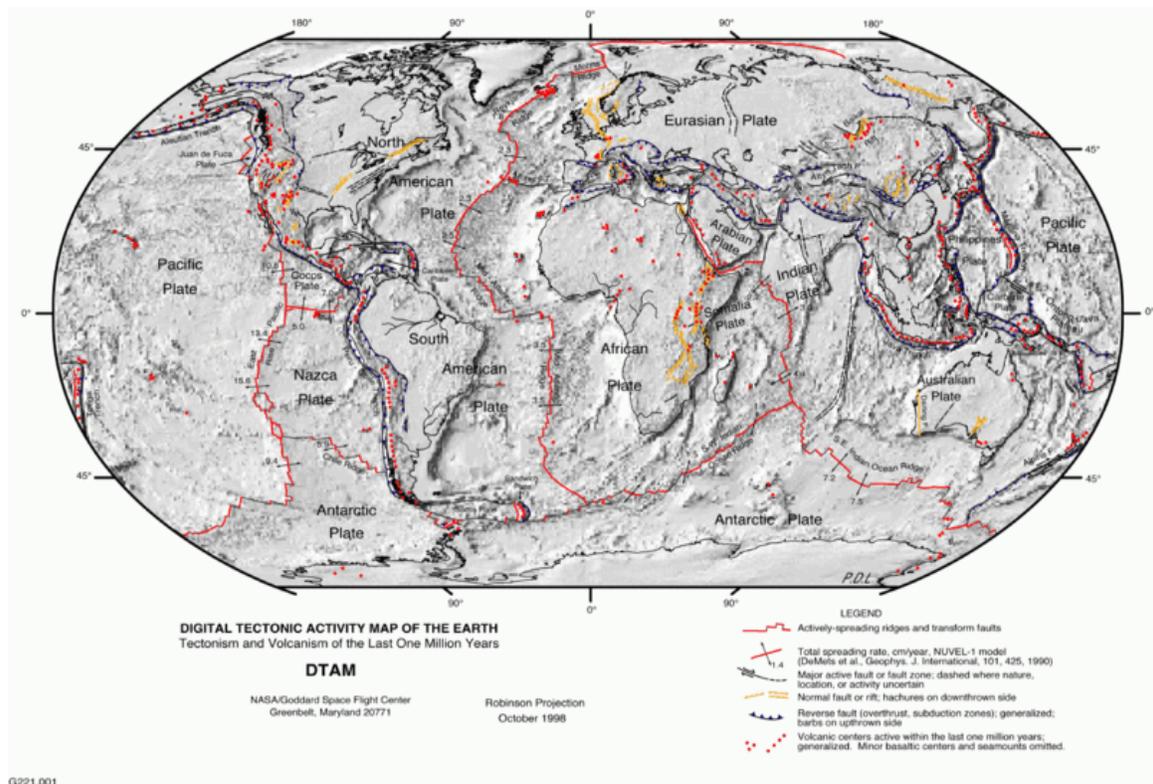
## Defining plate boundaries

Current plate boundaries are defined by their seismicity (Condie 1997). Past plate boundaries within existing plates are identified from evidence of vanished oceans, such as ophiolites (Lliboutry 2000).

## Past plate motions

The movement of plates has caused the formation and break-up of continents over time, including occasional formation of a supercontinent that contains most or all of the continents. The supercontinent Rodinia is thought to have formed about 1 billion years ago and to have embodied most or all of Earth's continents, and broken up into eight continents around 600 million years ago. The eight continents later re-assembled into another supercontinent called Pangaea; Pangaea broke up into Laurasia (which became North America and Eurasia) and Gondwana (which became the remaining continents).

Various types of quantitative and semi-quantitative information are available to constrain past plate motions. The geometric fit between continents, such as between west Africa and South America is still an important part of plate reconstruction. Magnetic stripe patterns provide a reliable guide to relative plate motions going back into the Jurassic period. The tracks of hotspots give absolute reconstructions but these are only available back to the Cretaceous (Torsvik 2008). Older reconstructions rely mainly on paleomagnetic pole data, although these only constrain the latitude and rotation, but not the longitude. Combining poles of different ages in a particular plate to produce apparent polar wander paths provides a method for comparing the motions of different plates through time (Butler 1992). Additional evidence comes from the distribution of certain sedimentary rock types, faunal provinces shown by particular fossil groups, and the position of orogenic belts (Torsvik 2008).

# Current plates



**DIGITAL TECTONIC ACTIVITY MAP OF THE EARTH**
Tectonism and Volcanism of the Last One Million Years

**DTAM**

NASA/Goddard Space Flight Center
Greenbelt, Maryland 20771

Robinson Projection
October 1998

LEGEND

G221.001

## Major plates

Depending on how they are defined, there are usually seven or eight "major" plates:

- African Plate
- Antarctic Plate
- Indo-Australian Plate, sometimes subdivided into:
  - Indian Plate
  - Australian Plate
- Eurasian Plate
- North American Plate
- South American Plate
- Pacific Plate

## Minor plates

There are dozens of smaller plates, the seven largest of which are:

- Arabian Plate
- Caribbean Plate

- Juan de Fuca Plate
- Cocos Plate
- Nazca Plate
- Philippine Sea Plate
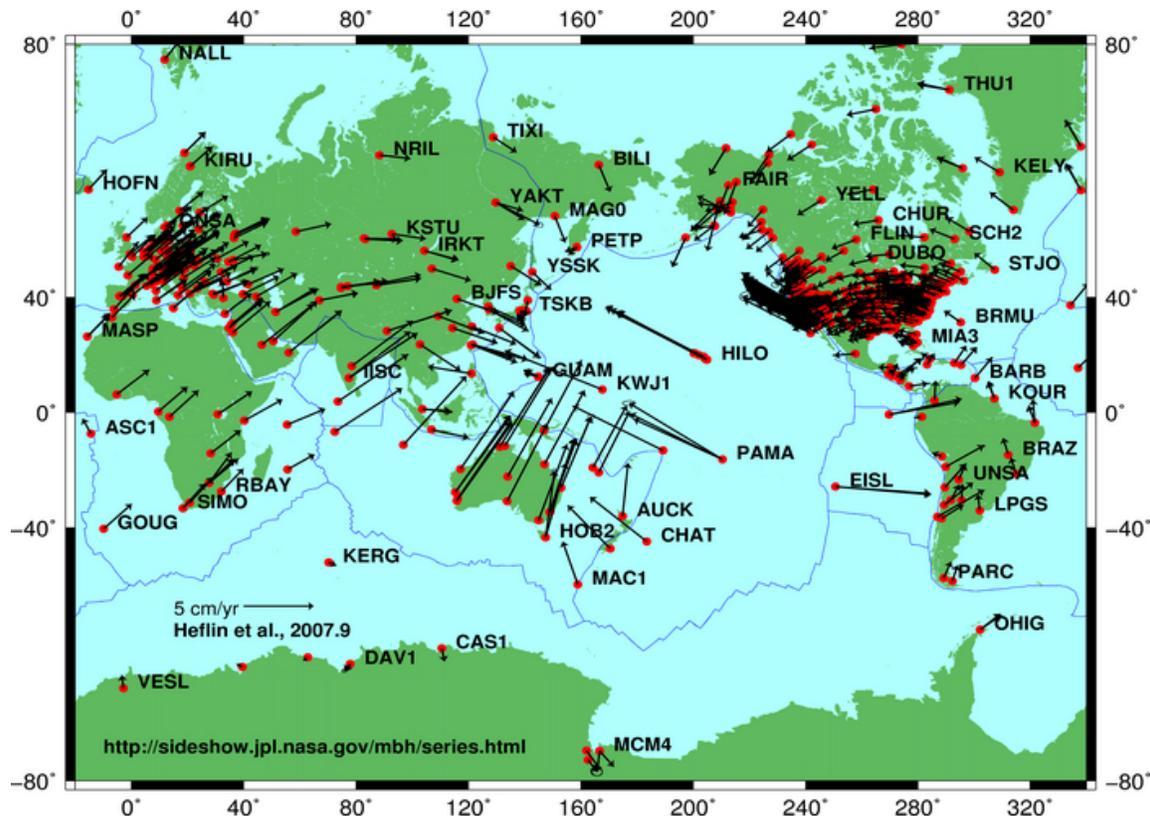- Scotia Plate

**Current Motion**

Plate motion based on Global Positioning System (GPS) satellite data from NASA JPL. The vectors show direction and magnitude of motion.

The current motion of the tectonic plates is nowadays revealed from remote sensing satellite data sets, calibrated with ground station measurements.

# Plate tectonics on other celestial bodies (Planets, Moons)

The appearance of plate tectonics on terrestrial planets is related to planetary mass, with more massive planets than Earth expected to exhibit plate tectonics. Earth may be a borderline case, owing its tectonic activity to abundant water (Valencia, O'Connell & Sasselov 2007)

(Silica and water form a deep eutectic.)

## Venus

Venus shows no evidence of active plate tectonics. There is debatable evidence of active tectonics in the planet's distant past; however, events taking place since then (such as the plausible and generally accepted hypothesis that the Venusian lithosphere has thickened greatly over the course of several hundred million years) has made constraining the course of its geologic record difficult. However, the numerous well-preserved impact craters have been utilized as a dating method to approximately date the Venusian surface (since there are thus far no known samples of Venusian rock to be dated by more reliable methods). Dates derived are dominantly in the range c. 500 to 750 Ma, although ages of up to c. 1.2 Ga have been calculated. This research has led to the fairly well accepted hypothesis that Venus has undergone an essentially complete volcanic resurfacing at least once in its distant past, with the last event taking place approximately within the range of estimated surface ages. While the mechanism of such an impressive thermal event remains a debated issue in Venusian geosciences, some scientists are advocates of processes involving plate motion to some extent.

One explanation for Venus' lack of plate tectonics is that on Venus temperatures are too high for significant water to be present (Kasting 1988). The Earth's crust is soaked with water, and water plays an important role in the development of shear zones. Plate tectonics requires weak surfaces in the crust along which crustal slices can move, and it may well be that such weakening never took place on Venus because of the absence of water. However, some researchers remain convinced that plate tectonics is or was once active on this planet.

## Mars

Mars is considerably smaller than Earth and Venus, and there is evidence for ice on its surface and in its crust.

In the 1990s, it was proposed that Martian Crustal Dichotomy was created by plate tectonic processes (Sleep 1994). Scientists today disagree, and believe that it was created either by upwelling within the Martian mantle that thickened the crust of the Southern Highlands and formed Tharsis (Zhong & Zuber 2001) or by a giant impact that excavated the Northern Lowlands (Andrews-Hanna, Zuber & Banerdt 2008).

Observations made of the magnetic field of Mars by the *Mars Global Surveyor* spacecraft in 1999 showed patterns of magnetic striping discovered on this planet. Some scientists interpreted these as requiring plate tectonic processes, such as seafloor spreading (Connerney et al. 1999, Connerney et al. 2005)). However, their data fail a "magnetic reversal test", which is used to see if they were formed by flipping polarities of a global magnetic field (Harrison 2000).

### Galilean satellites of Jupiter

Some of the satellites of Jupiter have features that may be related to plate-tectonic style deformation, although the materials and specific mechanisms may be different from plate-tectonic activity on Earth.
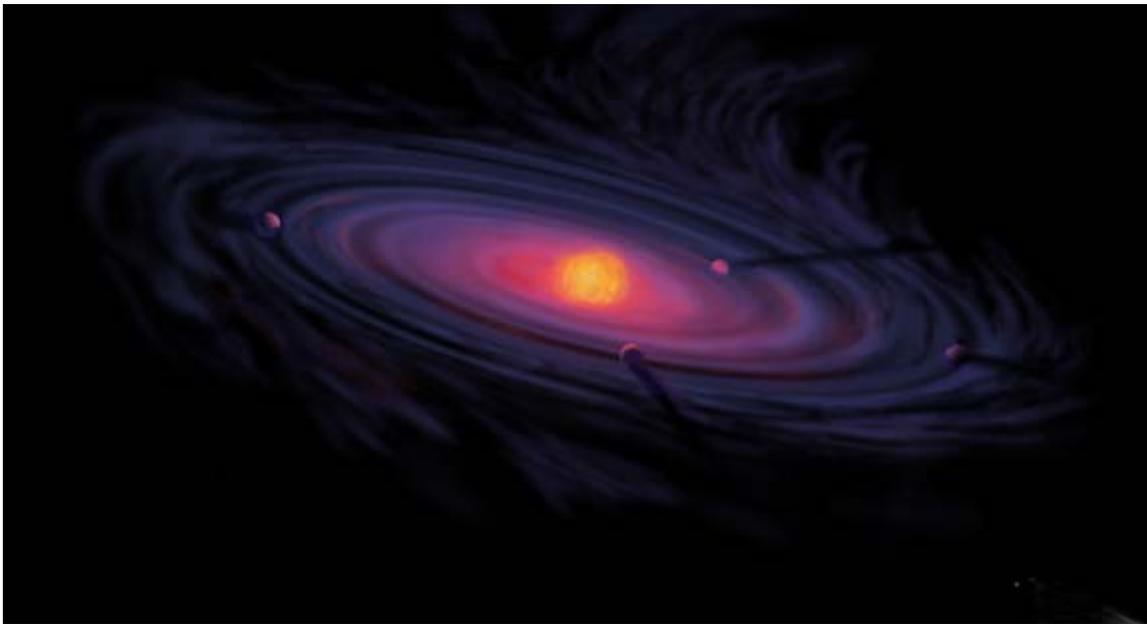
### Titan, moon of Saturn

Titan, the largest moon of Saturn, was reported to show tectonic activity in images taken by the Huygens Probe, which landed on Titan on January 14, 2005 (Soderblom et al. 2007).

### Exoplanets

It is believed that many planets around other stars will have plate tectonics. On Earth-sized planets, plate tectonics is more likely if there are oceans of water, but on larger super-earths plate tectonics is very likely even if the planet is dry (Valencia, O'Connell & Sasselov 2007).

# Chapter-4

# Formation and Evolution of the Solar System



Artist's conception of a protoplanetary disk

The formation and evolution of the Solar System is estimated to have begun 4.568 billion years ago with the gravitational collapse of a small part of a giant molecular cloud.

Most of the collapsing mass collected in the centre, forming the Sun, while the rest flattened into a protoplanetary disk out of which the planets, moons, asteroids, and other small Solar System bodies formed.

This widely accepted model, known as the nebular hypothesis, was first developed in the 18th century by Emanuel Swedenborg, Immanuel Kant, and Pierre-Simon Laplace. Its subsequent development has interwoven a variety of scientific disciplines including astronomy, physics, geology, and planetary science. Since the dawn of the space age in the 1950s and the discovery of extrasolar planets in the 1990s, the models have been both challenged and refined to account for new observations.

The Solar System has evolved considerably since its initial formation. Many moons have formed from circling discs of gas and dust around their parent planets, while other moons are believed to have formed independently and later been captured by their planets. Still others, as the Earth's Moon, may be the result of giant collisions. Collisions between bodies have occurred continually up to the present day and have been central to the evolution of the Solar System. The positions of the planets often shifted, and planets have switched places. This planetary migration now is believed to have been responsible for much of the Solar System's early evolution.

In roughly 5 billion years, the Sun will cool and expand outward to many times its current diameter (becoming a red giant), before casting off its outer layers as a planetary nebula, and leaving behind a stellar corpse known as a white dwarf. In the far distant future, the gravity of passing stars gradually will whittle away at the Sun's retinue of planets. Some planets will be destroyed, others ejected into interstellar space. Ultimately, over the course of trillions of years, it is likely that the Sun will be left alone with no bodies in orbit around it.

## History



Pierre-Simon Laplace, one of the originators of the nebular hypothesis

Ideas concerning the origin and fate of the world date from the earliest known writings; however, for almost all of that time, there was no attempt to link such theories to the existence of a "Solar System", simply because it was not generally believed that the Solar System, in the sense we now understand it, existed. The first step toward a theory of Solar System formation and evolution was the general acceptance of heliocentrism, which placed the Sun at the centre of the system and the Earth in orbit around it. This conception had gestated for millennia (philosophers such as Aristarchus of Samos had suggested it as early as 600 BC), but was widely accepted only by the end of the 17th century. The first recorded use of the term "Solar System" dates from 1704.

The current standard theory for Solar System formation, the nebular hypothesis, has fallen into and out of favour since its formulation by Emanuel Swedenborg, Immanuel Kant, and Pierre-Simon Laplace in the 18th century. The most significant criticism of the hypothesis was its apparent inability to explain the Sun's relative lack of angular momentum when compared to the planets. However, since the early 1980s studies of young stars have shown them to be surrounded by cool discs of dust and gas, exactly as the nebular hypothesis predicts, which has led to its re-acceptance.

Understanding of how the Sun will continue to evolve required an understanding of the source of its power. Arthur Stanley Eddington's confirmation of Albert Einstein's theory of relativity led to his realisation that the Sun's energy comes from nuclear fusion reactions in its core. In 1935, Eddington went further and suggested that other elements also might form within stars. Fred Hoyle elaborated on this premise by arguing that evolved stars called red giants created many elements heavier than hydrogen and helium in their cores. When a red giant finally casts off its outer layers, these elements would then be recycled to form other star systems.

# Formation

### Pre-solar nebula

The nebular hypothesis maintains that the Solar System formed from the gravitational collapse of a fragment of a giant molecular cloud. The cloud itself had a size of about 20 pc, while the fragments were roughly 1 pc (several light-years) across. The further collapse of the fragments led to the formation of dense cores 0.01–0.1 pc (2,000–20,000 AU) in size. One of these collapsing fragments (known as the *pre-solar nebula*) would form what became the Solar System. The composition of this region with a mass just over that of the Sun was about the same as that of the Sun today, with hydrogen, along with helium and trace amounts of lithium produced by Big Bang nucleosynthesis, forming about 98% of its mass. The remaining 2% of the mass consisted of heavier elements that were created by nucleosynthesis in earlier generations of stars. Late in the life of these stars, they ejected heavier elements into the interstellar medium.

Hubble image of protoplanetary discs in the Orion Nebula, a light-years-wide "stellar nursery" probably very similar to the primordial nebula from which our Sun formed

Studies of ancient meteorites reveal traces of stable daughter nuclei of short-lived isotopes, such as iron-60, that only form in exploding, short-lived stars. This indicates that one or more supernovae occurred near the Sun while it was forming. A shock wave from a supernova may have triggered the formation of the Sun by creating regions of over-density within the cloud, causing these regions to collapse. Because only massive, short-lived stars produce supernovae, the Sun must have formed in a large star-forming region that produced massive stars, possibly similar to the Orion Nebula. Studies of the structure of the Kuiper belt and of anomalous materials within it suggest that the Sun formed within a cluster of stars with a diameter of between 6.5 and 19.5 light-years and a collective mass equivalent to 3,000 Suns. Several simulations of our young Sun interacting with close-passing stars over the first 100 million years of its life produce anomalous orbits observed in the outer Solar System, such as detached objects.

Because of the conservation of angular momentum, the nebula spun faster as it collapsed. As the material within the nebula condensed, the atoms within it began to collide with increasing frequency, converting their kinetic energy into heat. The centre, where most of the mass collected, became increasingly hotter than the surrounding disc. Over about 100,000 years, the competing forces of gravity, gas pressure, magnetic fields, and rotation caused the contracting nebula to flatten into a spinning protoplanetary disc with a diameter of ~200 AU and form a hot, dense protostar (a star in which hydrogen fusion has not yet begun) at the centre.

At this point in its evolution, the Sun is believed to have been a T Tauri star. Studies of T Tauri stars show that they are often accompanied by discs of pre-planetary matter with masses of 0.001–0.1 solar masses. These discs extend to several hundred AU—the Hubble Space Telescope has observed protoplanetary discs of up to 1000 AU in diameter in star-forming regions such as the Orion Nebula—and are rather cool, reaching only a thousand kelvins at their hottest. Within 50 million years, the temperature and pressure at the core of the Sun became so great that its hydrogen began to fuse, creating an internal source of energy that countered gravitational contraction until hydrostatic equilibrium was achieved. This marked the Sun's entry into the prime phase of its life, known as the main sequence. Main sequence stars derive energy from the fusion of hydrogen into helium in their cores. The Sun remains a main sequence star today.

## Formation of planets



Artist's conception of the solar nebula

The various planets are thought to have formed from the *solar nebula*, the disc-shaped cloud of gas and dust left over from the Sun's formation. The currently accepted method by which the planets formed is known as accretion, in which the planets began as dust grains in orbit around the central protostar. Through direct contact, these grains formed into clumps up to 200 metres in diameter, which in turn collided to form larger bodies (planetesimals) of ~10 kilometres (km) in size. These gradually increased through further collisions, growing at the rate of centimetres per year over the course of the next few million years.

The inner Solar System, the region of the Solar System inside 4 AU, was too warm for volatile molecules like water and methane to condense, so the planetesimals that formed there could only form from compounds with high melting points, such as metals (like

iron, nickel, and aluminium) and rocky silicates. These rocky bodies would become the terrestrial planets (Mercury, Venus, Earth, and Mars). These compounds are quite rare in the universe, comprising only 0.6% of the mass of the nebula, so the terrestrial planets could not grow very large. The terrestrial embryos grew to about 0.05 Earth masses and ceased accumulating matter about 100,000 years after the formation of the Sun; subsequent collisions and mergers between these planet-sized bodies allowed terrestrial planets to grow to their present sizes.

When the terrestrial planets were forming, they remained immersed in a disk of gas and dust. The gas was partially supported by pressure and so did not orbit the Sun as rapidly as the planets. The resulting drag caused a transfer of angular momentum, and as a result the planets gradually migrated to new orbits. Models show that temperature variations in the disk governed this rate of migration, but the net trend was for the inner planets to migrate inward as the disk dissipated, leaving the planets in their current orbits.

The gas giant planets (Jupiter, Saturn, Uranus, and Neptune) formed further out, beyond the frost line, the point between the orbits of Mars and Jupiter where the material is cool enough for volatile icy compounds to remain solid. The ices that formed the Jovian planets were more abundant than the metals and silicates that formed the terrestrial planets, allowing the Jovian planets to grow massive enough to capture hydrogen and helium, the lightest and most abundant elements. Planetesimals beyond the frost line accumulated up to four Earth masses within about 3 million years. Today, the four gas giants comprise just under 99% of all the mass orbiting the Sun. Theorists believe it is no accident that Jupiter lies just beyond the frost line. Because the frost line accumulated large amounts of water via evaporation from infalling icy material, it created a region of lower pressure that increased the speed of orbiting dust particles and halted their motion toward the Sun. In effect, the frost line acted as a barrier that caused material to accumulate rapidly at ~5 AU from the Sun. This excess material coalesced into a large embryo of about 10 Earth masses, which then began to grow rapidly by swallowing hydrogen from the surrounding disc, reaching 150 Earth masses in only another 1000 years and finally topping out at 318 Earth masses. Saturn may owe its substantially lower mass simply to having formed a few million years after Jupiter, when there was less gas available to consume.

T Tauri stars like the young Sun have far stronger stellar winds than more stable, older stars. Uranus and Neptune are believed to have formed after Jupiter and Saturn did, when the strong solar wind had blown away much of the disc material. As a result, the planets accumulated little hydrogen and helium—not more than 1 Earth mass each. Uranus and Neptune are sometimes referred to as failed cores. The main problem with formation theories for these planets is the timescale of their formation. At the current locations it would have taken a hundred million years for their cores to accrete. This means that Uranus and Neptune probably formed closer to the Sun—near or even between Jupiter and Saturn—and later migrated outward. Motion in the planetesimal era was not all inward toward the Sun; the *Stardust* sample return from Comet Wild 2 has suggested that materials from the early formation of the Solar System migrated from the warmer inner Solar System to the region of the Kuiper belt.

After between three and ten million years, the young Sun's solar wind would have cleared away all the gas and dust in the protoplanetary disc, blowing it into interstellar space, thus ending the growth of the planets.

# Subsequent evolution



The giant impact believed to have formed the Moon.

The planets were originally believed to have formed in or near their current orbits. However, this view underwent radical change during the late 20th and early 21st centuries. Currently, it is believed that the Solar System looked very different after its initial formation: several objects at least as massive as Mercury were present in the inner Solar System, the outer Solar System was much more compact than it is now, and the Kuiper belt was much closer to the Sun.

### Terrestrial planets

At the end of the planetary formation epoch the inner Solar System was populated by 50–100 Moon- to Mars-sized planetary embryos. Further growth was possible only because these bodies collided and merged, which took less than 100 million years. These objects would have gravitationally interacted with one another, tugging at each other's orbits until they collided, growing larger until the four terrestrial planets we know today took shape. One such giant collision is believed to have formed the Moon, while another removed the outer envelope of the young Mercury.

One unresolved issue with this model is that it cannot explain how the initial orbits of the proto-terrestrial planets, which would have needed to be highly eccentric to collide, produced the remarkably stable and near-circular orbits the terrestrial planets possess today. One hypothesis for this "eccentricity dumping" is that the terrestrials formed in a disc of gas still not expelled by the Sun. The "gravitational drag" of this residual gas would have eventually lowered the planets' energy, smoothing out their orbits. However,

such gas, if it existed, would have prevented the terrestrials' orbits from becoming so eccentric in the first place. Another hypothesis is that gravitational drag occurred not between the planets and residual gas but between the planets and the remaining small bodies. As the large bodies moved through the crowd of smaller objects, the smaller objects, attracted by the larger planets' gravity, formed a region of higher density, a "gravitational wake", in the larger objects' path. As they did so, the increased gravity of the wake slowed the larger objects down into more regular orbits.
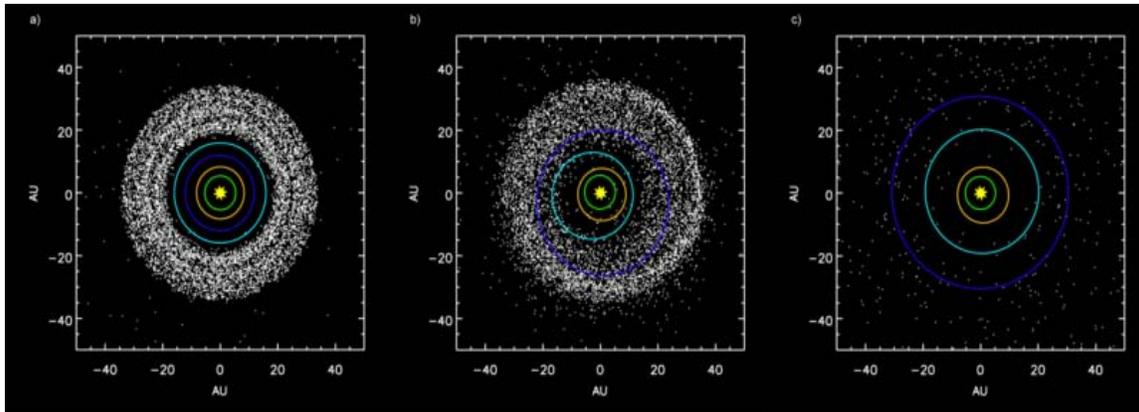
## Asteroid belt

The outer edge of the terrestrial region, between 2 and 4 AU from Sun, is called the asteroid belt. The asteroid belt initially contained more than enough matter to form 2–3 Earth-like planets, and, indeed, a large number of planetesimals formed there. As with the terrestrials, planetesimals in this region later coalesced and formed 20–30 Moon- to Mars-sized planetary embryos; however, the proximity of Jupiter meant that after this planet formed, 3 million years after the Sun, the region's history changed dramatically. Orbital resonances with Jupiter and Saturn are particularly strong in the asteroid belt, and gravitational interactions with more massive embryos scattered many planetesimals into those resonances. Jupiter's gravity increased the velocity of objects within these resonances, causing them to shatter upon collision with other bodies, rather than accrete.

As Jupiter migrated inward following its formation, resonances would have swept across the asteroid belt, dynamically exciting the region's population and increasing their velocities relative to each other. The cumulative action of the resonances and the embryos either scattered the planetesimals away from the asteroid belt or excited their orbital inclinations and eccentricities. Some of those massive embryos too were ejected by Jupiter, while others may have migrated to the inner Solar System and played a role in the final accretion of the terrestrial planets. During this primary depletion period, the effects of the giant planets and planetary embryos left the asteroid belt with a total mass equivalent to less than 1% that of the Earth, composed mainly of small planetesimals. This is still 10–20 times more than the current mass in the main belt, which is about 1/2,000 the Earth's mass. A secondary depletion period that brought the asteroid belt down close to its present mass is believed to have followed when Jupiter and Saturn entered a temporary 2:1 orbital resonance (see below).

The inner Solar System's period of giant impacts probably played a role in the Earth acquiring its current water content ($\sim 6 \times 10^{21}$ kg) from the early asteroid belt. Water is too volatile to have been present at Earth's formation and must have been subsequently delivered from outer, colder parts of the Solar System. The water was probably delivered by planetary embryos and small planetesimals thrown out of the asteroid belt by Jupiter. A population of main-belt comets discovered in 2006 has been also suggested as a possible source for Earth's water. In contrast, comets from the Kuiper belt or farther regions delivered not more than about 6% of Earth's water. The panspermia hypothesis holds that life itself may have been deposited on Earth in this way, although this idea is not widely accepted.

## Planetary migration

According to the nebular hypothesis, the outer two planets are in the "wrong place". Uranus and Neptune (known as the "ice giants") exist in a region where the reduced density of the solar nebula and longer orbital times render their formation highly implausible. The two are instead believed to have formed in orbits near Jupiter and Saturn, where more material was available, but to have migrated outward to their current positions over hundreds of millions of years.



Simulation showing outer planets and Kuiper belt: a) Before Jupiter/Saturn 2:1 resonance b) Scattering of Kuiper belt objects into the Solar System after the orbital shift of Neptune c) After ejection of Kuiper belt bodies by Jupiter

The migration of the outer planets is also necessary to account for the existence and properties of the Solar System's outermost regions. Beyond Neptune, the Solar System continues into the Kuiper belt, the scattered disc, and the Oort cloud, three sparse populations of small icy bodies thought to be the points of origin for most observed comets. At their distance from the Sun, accretion was too slow to allow planets to form before the solar nebula dispersed, and thus the initial disc lacked enough mass density to consolidate into a planet. The Kuiper belt lies between 30 and 55 AU from the Sun, while the farther scattered disc extends to over 100 AU, and the distant Oort cloud begins at about 50,000 AU. Originally, however, the Kuiper belt was much denser and closer to the Sun, with an outer edge at approximately 30 AU. Its inner edge would have been just beyond the orbits of Uranus and Neptune, which were in turn far closer to the Sun when they formed (most likely in the range of 15–20 AU), and in opposite locations, with Uranus farther from the Sun than Neptune.

After the formation of the Solar System, the orbits of all the giant planets continued to change slowly, influenced by their interaction with large number of remaining planetesimals. After 500–600 million years (about 4 billion years ago) Jupiter and Saturn fell into a 2:1 resonance; Saturn orbited the Sun once for every two Jupiter orbits. This resonance created a gravitational push against the outer planets, causing Neptune to surge past Uranus and plough into the ancient Kuiper belt. The planets scattered the majority of

the small icy bodies inwards, while themselves moving outwards. These planetesimals then scattered off the next planet they encountered in a similar manner, moving the planets' orbits outwards while they moved inwards. This process continued until the planetesimals interacted with Jupiter, whose immense gravity sent them into highly elliptical orbits or even ejected them outright from the Solar System. This caused Jupiter to move slightly inward. Those objects scattered by Jupiter into highly elliptical orbits formed the Oort cloud; those objects scattered to a lesser degree by the migrating Neptune formed the current Kuiper belt and scattered disc. This scenario explains the Kuiper belt's and scattered disc's present low mass. Some of the scattered objects, including Pluto, became gravitationally tied to Neptune's orbit, forcing them into mean-motion resonances. Eventually, friction within the planetesimal disc made the orbits of Uranus and Neptune circular again.

In contrast to the outer planets, the inner planets are not believed to have migrated significantly over the age of the Solar System, because their orbits have remained stable following the period of giant impacts.

## Late Heavy Bombardment and after



Meteor Crater in Arizona. Created 50,000 years ago by an impactor only 50m across, it is a stark reminder that the accretion of the Solar System is not over.

Gravitational disruption from the outer planets' migration would have sent large numbers of asteroids into the inner Solar System, severely depleting the original belt until it reached today's extremely low mass. This event may have triggered the Late Heavy

Bombardment that occurred approximately 4 billion years ago, 500–600 million years after the formation of the Solar System. This period of heavy bombardment lasted several hundred million years and is evident in the cratering still visible on geologically dead bodies of the inner Solar System such as the Moon and Mercury. The oldest known evidence for life on Earth dates to 3.8 billion years ago—almost immediately after the end of the Late Heavy Bombardment.

Impacts are believed to be a regular (if currently infrequent) part of the evolution of the Solar System. That they continue to happen is evidenced by the collision of Comet Shoemaker-Levy 9 with Jupiter in 1994, the 2009 Jupiter impact event, and the impact feature Meteor Crater in Arizona. The process of accretion, therefore, is not complete, and may still pose a threat to life on Earth.

Over the course of the Solar System's evolution, comets were ejected out of the inner Solar System by the gravity of the giant planets, and sent thousands of AU outward to form the Oort cloud, a spherical outer swarm of cometary nuclei at the farthest extent of the Sun's gravitational pull. Eventually, after about 800 million years, the gravitational disruption caused by galactic tides, passing stars and giant molecular clouds began to deplete the cloud, sending comets into the inner Solar System. The evolution of the outer Solar System also appears to have been influenced by space weathering from the solar wind, micrometeorites, and the neutral components of the interstellar medium.

The evolution of the asteroid belt after Late Heavy Bombardment was mainly governed by collisions. Objects with large mass have enough gravity to retain any material ejected by a violent collision. In the asteroid belt this usually is not the case. As a result, many larger objects have been broken apart, and sometimes newer objects have been forged from the remnants in less violent collisions. Moons around some asteroids currently can only be explained as consolidations of material flung away from the parent object without enough energy to entirely escape its gravity.

# Moons

Moons have come to exist around most planets and many other Solar System bodies. These natural satellites originated by one of three possible mechanisms:

- co-formation from a circum-planetary disc (only in the cases of the gas giants);
- formation from impact debris (given a large enough impact at a shallow angle); and
- capture of a passing object.

Jupiter and Saturn have a number of large moons, such as Io, Europa, Ganymede and Titan, which may have originated from discs around each giant planet in much the same way that the planets formed from the disc around the Sun. This origin is indicated by the large sizes of the moons and their proximity to the planet. These attributes are impossible to achieve via capture, while the gaseous nature of the primaries make formation from collision debris another impossibility. The outer moons of the gas giants tend to be small

and have eccentric orbits with arbitrary inclinations. These are the characteristics expected of captured bodies. Most such moons orbit in the direction opposite the rotation of their primary. The largest irregular moon is Neptune's moon Triton, which is believed to be a captured Kuiper belt object.

Moons of solid Solar System bodies have been created by both collisions and capture. Mars's two small moons, Deimos and Phobos, are believed to be captured asteroids. The Earth's Moon is believed to have formed as a result of a single, large oblique collision. The impacting object likely had a mass comparable to that of Mars, and the impact probably occurred near the end of the period of giant impacts. The collision kicked into orbit some of the impactor's mantle, which then coalesced into the Moon. The impact was probably the last in the series of mergers that formed the Earth. It has been further hypothesized that the Mars-sized object may have formed at one of the stable Earth-Sun Lagrangian points (either $L_4$ or $L_5$) and drifted from its position. Pluto's moon Charon may also have formed by means of a large collision; the Pluto-Charon and Earth-Moon systems are the only two in the Solar System in which the satellite's mass is at least 1% that of the larger body.

# Future

Astronomers estimate that the Solar System as we know it today will not change drastically until the Sun has fused all the hydrogen fuel in its core into helium, beginning its evolution from the main sequence of the Hertzsprung-Russell diagram and into its red giant phase. Even so, the Solar System will continue to evolve until then.

### Long-term stability

The Solar System is chaotic, with the orbits of the planets open to long-term variations. One notable example of this chaos is the Neptune-Pluto system, which lies in a 3:2 orbital resonance. Although the resonance itself will remain stable, it becomes impossible to predict the position of Pluto with any degree of accuracy more than 10–20 million years (the Lyapunov time) into the future. Another example is Earth's axial tilt which, thanks to friction raised within Earth's mantle by tidal interactions with the Moon (see below) will be incomputable at some point between 1.5 and 4.5 billion years from now.

The outer planets' orbits are chaotic over longer timescales, such that they possess a Lyapunov time in the range of 2–230 million years. In all cases this means that the position of a planet along its orbit ultimately becomes impossible to predict with any certainty (so, for example, the timing of winter and summer become uncertain), but in some cases the orbits themselves may change dramatically. Such chaos manifests most strongly as changes in eccentricity, with some planets' orbits becoming significantly more—or less—elliptical.

Ultimately, the Solar System is stable in that none of the planets will collide with each other or be ejected from the system in the next few billion years. Beyond this, within five billion years or so Mars's eccentricity may grow to around 0.2, such that it lies on an

Earth-crossing orbit, leading to a potential collision. In the same timescale, Mercury's eccentricity may grow even further, and a close encounter with Venus could theoretically eject it from the Solar System altogether or send it on a collision course with Venus or Earth.

## Moon-ring systems

The evolution of moon systems is driven by tidal forces. A moon will raise a tidal bulge in the object it orbits (the primary) due to the differential gravitational force across diameter of the primary. If a moon is revolving in the same direction as the planet's rotation and the planet is rotating faster than the orbital period of the moon, the bulge will constantly be pulled ahead of the moon. In this situation, angular momentum is transferred from the rotation of the primary to the revolution of the satellite. The moon gains energy and gradually spirals outward, while the primary rotates more slowly over time.

The Earth and its Moon are one example of this configuration. Today, the Moon is tidally locked to the Earth; one of its revolutions around the Earth (currently about 29 days) is equal to one of its rotations about its axis, so it always shows one face to the Earth. The Moon will continue to recede from Earth, and Earth's spin will continue to slow gradually. In about 50 billion years, if they survive the Sun's expansion, the Earth and Moon will become tidally locked to each other; each will be caught up in what is called a "spin–orbit resonance" in which the Moon will circle the Earth in about 47 days and both Moon and Earth will rotate around their axes in the same time, each only visible from one hemisphere of the other. Other examples are the Galilean moons of Jupiter (as well as many of Jupiter's smaller moons) and most of the larger moons of Saturn.

Neptune and its moon Triton, taken by *Voyager 2*. Triton's orbit will eventually take it within Neptune's Roche limit, tearing it apart and possibly forming a new ring system.
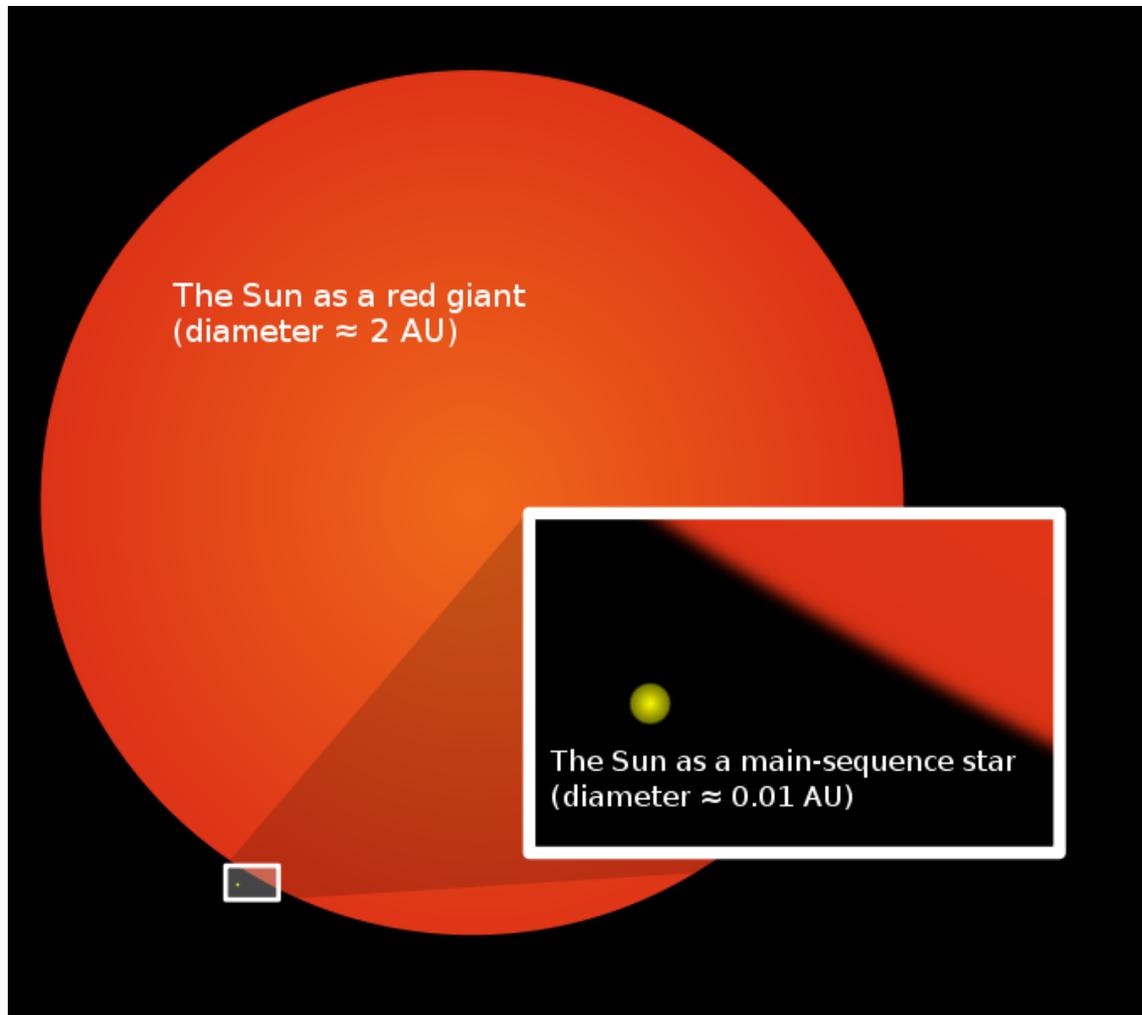
A different scenario occurs when the moon is either revolving around the primary faster than the primary rotates, or is revolving in the direction opposite the planet's rotation. In these cases, the tidal bulge lags behind the moon in its orbit. In the former case, the direction of angular momentum transfer is reversed, so the rotation of the primary speeds up while the satellite's orbit shrinks. In the latter case, the angular momentum of the rotation and revolution have opposite signs, so transfer leads to decreases in the magnitude of each (that cancel each other out). In both cases, tidal deceleration causes the moon to spiral in towards the primary until it either is torn apart by tidal stresses, potentially creating a planetary ring system, or crashes into the planet's surface or atmosphere. Such a fate awaits the moons Phobos of Mars (within 30 to 50 million years), Triton of Neptune (in 3.6 billion years), Metis and Adrastea of Jupiter, and at least 16 small satellites of Uranus and Neptune. Uranus' Desdemona may even collide with one of its neighboring moons.

A third possibility is where the primary and moon are tidally locked to each other. In that case, the tidal bulge stays directly under the moon, there is no transfer of angular momentum, and the orbital period will not change. Pluto and Charon are an example of this type of configuration.

Prior to the 2004 arrival of the *Cassini–Huygens* spacecraft, the rings of Saturn were widely thought to be much younger than the Solar System and were not expected to survive beyond another 300 million years. Gravitational interactions with Saturn's moons were expected to gradually sweep the rings' outer edge toward the planet, with abrasion by meteorites and Saturn's gravity eventually taking the rest, leaving Saturn unadorned. However, data from the *Cassini* mission led scientists to revise that early view. Observations revealed 10 km-wide icy clumps of material that repeatedly break apart and reform, keeping the rings fresh. Saturn's rings are far more massive than the rings of the other gas giants. This large mass is believed to have preserved Saturn's rings since the planet first formed 4.5 billion years ago, and is likely to preserve them for billions of years to come.

## The Sun and planetary environments

In the long term, the greatest changes in the Solar System will come from changes in the Sun itself as it ages. As the Sun burns through its supply of hydrogen fuel, it gets hotter and burns the remaining fuel even faster. As a result, the Sun is growing brighter at a rate of ten percent every 1.1 billion years. In one billion years' time, as the Sun's radiation output increases, its circumstellar habitable zone will move outwards, making the Earth's surface hot enough that liquid water can no longer exist there naturally. At this point, all life on land will become extinct. Evaporation of water, a potent greenhouse gas, from the oceans' surface could accelerate temperature increase, potentially ending all life on Earth even sooner. During this time, it is possible that as Mars's surface temperature gradually rises, carbon dioxide and water currently frozen under the surface soil will release into the atmosphere, creating a greenhouse effect that will heat the planet until it achieves conditions parallel to Earth today, providing a potential future abode for life. By 3.5 billion years from now, Earth's surface conditions will be similar to those of Venus today.

The Sun as a red giant
(diameter ≈ 2 AU)
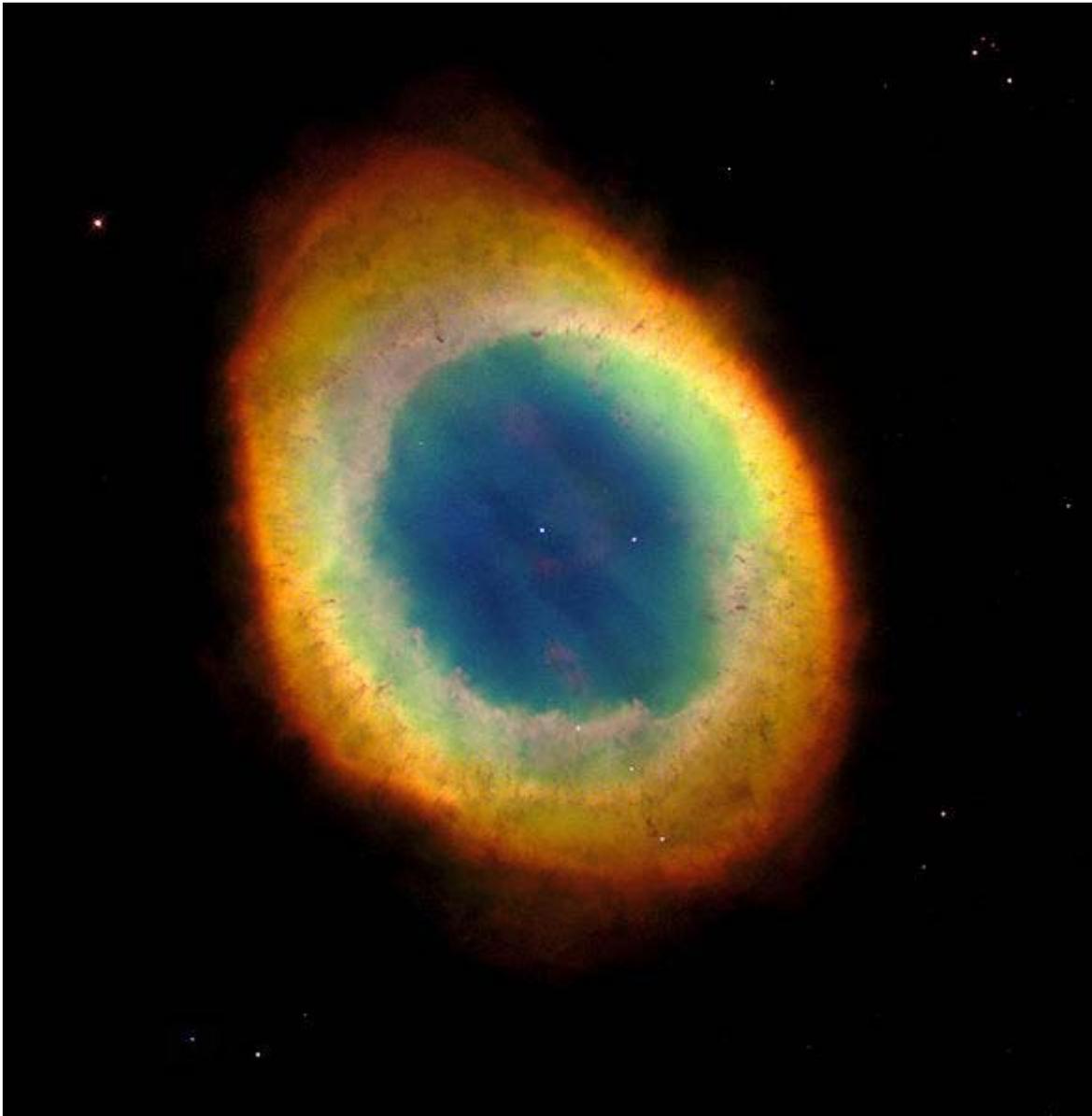
The Sun as a main-sequence star
(diameter ≈ 0.01 AU)

Relative size of our Sun as it is now (inset) compared to its estimated future size as a red giant

Around 5.4 billion years from now, the core of the Sun will become hot enough to trigger hydrogen fusion in its surrounding shell. This will cause the outer layers of the star to expand greatly, and the star will enter a phase of its life in which it is called a red giant. Within 7.5 billion years, the Sun will have expanded to a radius of 1.2 AU—256 times its current size. At the tip of the red giant branch, as a result of the vastly increased surface area, the Sun's surface will be much cooler (about 2600 K) than now and its luminosity much higher—up to 2,700 current solar luminosities. For part of its red giant life, the Sun will have a strong stellar wind that will carry away around 33% of its mass. During these times, it is possible that Saturn's moon Titan could achieve surface temperatures necessary to support life.

As the Sun expands, it will swallow the planets Mercury and, most likely, Venus. Earth's fate is less clear; although the Sun will envelop Earth's current orbit, the star's loss of mass (and thus weaker gravity) will cause the planets' orbits to move farther out. If it were only for this, Venus and Earth would probably escape incineration, but a 2008 study

suggests that Earth will likely be swallowed up as a result of tidal interactions with the Sun's weakly bound outer envelope.

Gradually, the hydrogen burning in the shell around the solar core will increase the mass of the core until it reaches about 45% of the present solar mass. At this point the density and temperature will become so high that the fusion of helium into carbon will begin, leading to a helium flash; the Sun will shrink from around 250 to 11 times its present (main sequence) radius. Consequently, its luminosity will decrease from around 3,000 to 54 times its current level, and its surface temperature will increase to about 4770 K. The Sun will become a horizontal branch star, burning helium in its core in a stable fashion much like it burns hydrogen today. The helium-fusing stage will last only 100 million years. Eventually, it will have to again resort to the reserves of hydrogen and helium in its outer layers and will expand a second time, turning into what is known as an asymptotic giant branch star. Here the luminosity of the Sun will increase again, reaching about 2,090 present luminosities, and it will cool to about 3500 K. This phase lasts about 30 million years, after which, over the course of a further 100,000 years, the Sun's remaining outer layers will fall away, ejecting a vast stream of matter into space and forming a halo known (misleadingly) as a planetary nebula. The ejected material will contain the helium and carbon produced by the Sun's nuclear reactions, continuing the enrichment of the interstellar medium with heavy elements for future generations of stars.
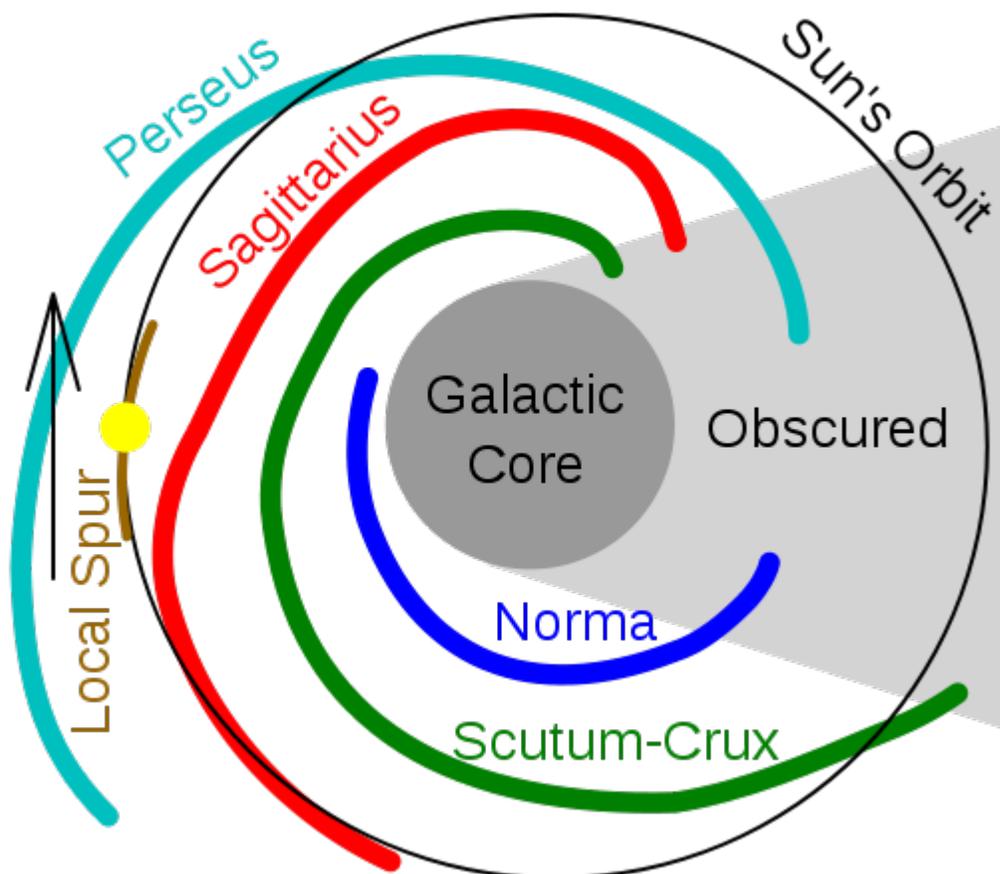
The Ring nebula, a planetary nebula similar to what the Sun will become

This is a relatively peaceful event, nothing akin to a supernova, which our Sun is too small to undergo as part of its evolution. Any observer present to witness this occurrence would see a massive increase in the speed of the solar wind, but not enough to destroy a planet completely. However, the star's loss of mass could send the orbits of the surviving planets into chaos, causing some to collide, others to be ejected from the Solar System, and still others to be torn apart by tidal interactions. Afterwards, all that will remain of the Sun is a white dwarf, an extraordinarily dense object, 54% its original mass but only the size of the Earth. Initially, this white dwarf may be 100 times as luminous as the Sun is now. It will consist entirely of degenerate carbon and oxygen, but will never reach

temperatures hot enough to fuse these elements. Thus the white dwarf Sun will gradually cool, growing dimmer and dimmer.

As the Sun dies, its gravitational pull on the orbiting bodies such as planets, comets and asteroids will weaken due to its mass loss. All remaining planets' orbits will expand; if Venus, Earth, and Mars still exist, their orbits will lie roughly at 1.4 AU (210,000,000 km), 1.9 AU (280,000,000 km), and 2.8 AU (420,000,000 km). They and the other remaining planets will become dark, frigid hulks, completely devoid of any form of life. They will continue to orbit their star, their speed slowed due to their increased distance from the Sun and the Sun's reduced gravity. Two billion years later, when the Sun has cooled to the 6000–8000K range, the carbon and oxygen in the Sun's core will freeze, with over 90% of its remaining mass assuming a crystalline structure. Eventually, after trillions more years, the Sun will finally cease to shine altogether, becoming a black dwarf.

## Galactic interaction



Location of the Solar System within our galaxy

The Solar System travels alone through the Milky Way galaxy in a circular orbit approximately 30,000 light years from the galactic centre. Its speed is about 220 km/s. The period required for the Solar System to complete one revolution around the galactic centre, the galactic year, is in the range of 220–250 million years. Since its formation, the Solar System has completed at least 20 such revolutions.

A number of scientists have speculated that the Solar System's path through the galaxy is a factor in the periodicity of mass extinctions observed in the Earth's fossil record. One hypothesis supposes that vertical oscillations made by the Sun as it orbits the galactic centre cause it to regularly pass through the galactic plane. When the Sun's orbit takes it outside the galactic disc, the influence of the galactic tide is weaker; as it re-enters the galactic disc, as it does every 20–25 million years, it comes under the influence of the far stronger "disc tides", which, according to mathematical models, increase the flux of Oort cloud comets into the Solar System by a factor of 4, leading to a massive increase in the likelihood of a devastating impact.

However, others argue that the Sun is currently close to the galactic plane, and yet the last great extinction event was 15 million years ago. Therefore the Sun's vertical position cannot alone explain such periodic extinctions, and that extinctions instead occur when the Sun passes through the galaxy's spiral arms. Spiral arms are home not only to larger numbers of molecular clouds, whose gravity may distort the Oort cloud, but also to higher concentrations of bright blue giant stars, which live for relatively short periods and then explode violently as supernovae.

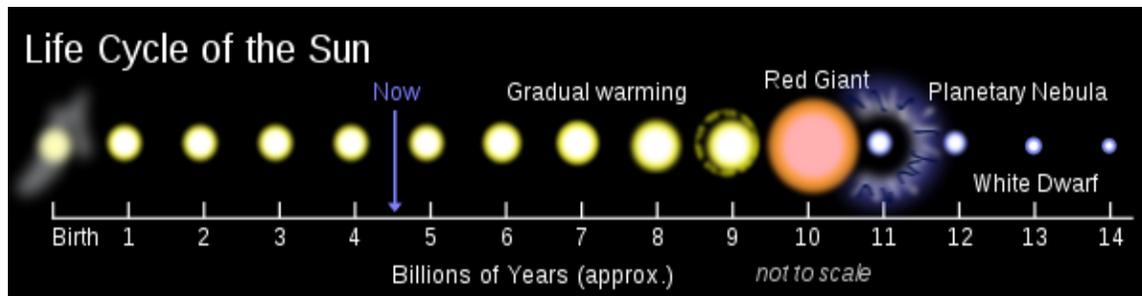## Galactic collision and planetary disruption

Although the vast majority of galaxies in the Universe are moving away from the Milky Way, the Andromeda Galaxy, the largest member of our Local Group of galaxies, is heading towards it at about 120 km/s. In 2 billion years, Andromeda and the Milky Way will collide, causing both to deform as tidal forces distort their outer arms into vast tidal tails. If this initial disruption occurs, astronomers calculate a 12% chance that the Solar System will be pulled outward into the Milky Way's tidal tail and a 3% chance that it will become gravitationally bound to Andromeda and thus a part of that galaxy. After a further series of glancing blows, during which the likelihood of the Solar System's ejection rises to 30%, the galaxies' supermassive black holes will merge. Eventually, in roughly 7 billion years, the Milky Way and Andromeda will complete their merger into a giant elliptical galaxy. During the merger, if there is enough gas, the increased gravity will force the gas to the centre of the forming elliptical galaxy. This may lead to a short period of intensive star formation called a starburst. In addition the infalling gas will feed the newly formed black hole transforming it into an active galactic nucleus. The force of these interactions will likely push the Solar System into the new galaxy's outer halo, leaving it relatively unscathed by the radiation from these collisions.

It is a common misconception that this collision will disrupt the orbits of the planets in the Solar System. While it is true that the gravity of passing stars can detach planets into interstellar space, distances between stars are so great that the likelihood of the Milky

Way-Andromeda collision causing such disruption to any individual star system is negligible. While the Solar System as a whole could be affected by these events, the Sun and planets are not expected to be disturbed.

However, over time, the cumulative probability of a chance encounter with a star increases, and disruption of the planets becomes all but inevitable. Assuming that the Big Crunch or Big Rip scenarios for the end of the universe do not occur, calculations suggest that the gravity of passing stars will have completely stripped the dead Sun of its remaining planets within 1 quadrillion ($10^{15}$) years. This point marks the end of the Solar System. While the Sun and planets may survive, the Solar System, in any meaningful sense, will cease to exist.

# Chronology



The time frame of the Solar System's formation has been determined using radiometric dating. Scientists estimate that the Solar System is 4.6 billion years old. The oldest known mineral grains on Earth are approximately 4.4 billion years old. Rocks this old are rare, as Earth's surface is constantly being reshaped by erosion, volcanism, and plate tectonics. To estimate the age of the Solar System, scientists use meteorites, which were formed during the early condensation of the solar nebula. Almost all meteorites are found to have an age of 4.6 billion years, suggesting that the Solar System must be at least this old.

Studies of discs around other stars have also done much to establish a time frame for Solar System formation. Stars between one and three million years old possess discs rich in gas, whereas discs around stars more than 10 million years old have little to no gas, suggesting that gas giant planets within them have ceased forming.

# Timeline of Solar System evolution

Note: All dates and times in this chronology are approximate and should be taken as an order of magnitude indicator only.

| Phase | Time since formation of the Sun | Event |
|---|---|---|
| **Pre-Solar System** | Billions of years before the formation of the Solar System | Previous generations of stars live and die, injecting heavy elements into the interstellar medium out of which the Solar System formed. |
| | ~ 50 million years before formation of the Solar System | If the Solar System formed in an Orion nebula-like star-forming region, the most massive stars are formed, live their lives, die, and explode in supernovae. One particular supernova, called the *primal supernova*, possibly triggers the formation of the Solar System. |
| **Formation of Sun** | **0**–100,000 years | Pre-solar nebula forms and begins to collapse. Sun begins to form. |
| | 100,000 – 50 million years | Sun is a T Tauri protostar. |
| | 100,000 - 10 million years | Outer planets form. By 10 million years, gas in the protoplanetary disc has been blown away, and outer planet formation is likely complete. |
| | 10 million - 100 million years | Terrestrial planets and the Moon form. Giant impacts occur. Water delivered to Earth. |
| **Main sequence** | 50 million years | Sun becomes a main sequence star. |
| | 200 million years | Oldest known rocks on the Earth formed. |
| | 500 million – 600 million years | Resonance in Jupiter and Saturn's orbits moves Neptune out into the Kuiper belt. Late Heavy Bombardment occurs in the inner Solar System. |
| | 800 million years | Oldest known life on Earth. Oort cloud reaches maximum mass. |
| | 4.6 billion years | **Today**. Sun remains a main sequence star, continually growing warmer and brighter by ~10% every 1 billion years. |
| | 6 billion years | Sun's habitable zone moves outside of the Earth's orbit, possibly shifting onto Mars' orbit. |
| | 7 billion years | The Milky Way and Andromeda Galaxy begin to collide. Slight chance the Solar System could be captured by Andromeda before the two galaxies fuse completely. |
| **Post-main** | 10 billion – 12 | Sun starts burning hydrogen in a shell surrounding its |

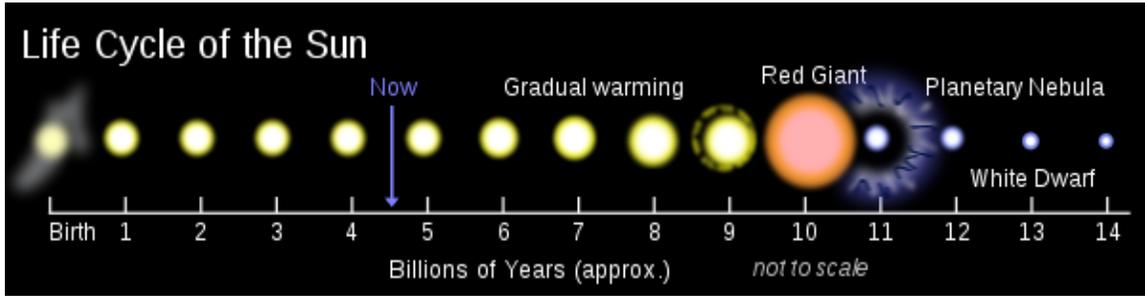| | | |
|---|---|---|
| sequence | billion years | core, ending its main sequence life. Sun begins to ascend the red giant branch of the Hertzsprung-Russell diagram, growing dramatically more luminous (by a factor of up to 2,700), larger (by a factor of up to 250 in radius), and cooler (down to 2600 K): Sun is now a red giant. Mercury and possibly Venus and Earth are swallowed. Saturn's moon Titan may become habitable. |
| | ~ 12 billion years | Sun passes through helium-burning horizontal branch and asymptotic giant branch phases, losing a total of ~30% of its mass in all post-main sequence phases. Asymptotic giant branch phase ends with the ejection of a planetary nebula, leaving the core of the Sun behind as a white dwarf. |
| **Remnant Sun** | > 12 billion years | The white dwarf Sun, no longer producing energy, begins to cool and dim continuously; this continues for trillions of years, eventually reaching a black dwarf state. |
| | ~ 1 quadrillion years ($10^{15}$ years) | Sun cools to 5 K. Gravity of passing stars detaches planets from orbits. Solar System ceases to exist. |

# Chapter-5

# Stellar Evolution



Life cycle of a Sun-like star.

**Stellar evolution** is the process by which a star undergoes a sequence of radical changes during its lifetime. Depending on the mass of the star, this lifetime ranges from only a few million years (for the most massive) to trillions of years (for the least massive, which is considerably more than the age of the universe).

Stellar evolution is not studied by observing the life of a single star, as most stellar changes occur too slowly to be detected, even over many centuries. Instead, astrophysicists come to understand how stars evolve by observing numerous stars at the various points in their life, and by simulating stellar structure with computer models.

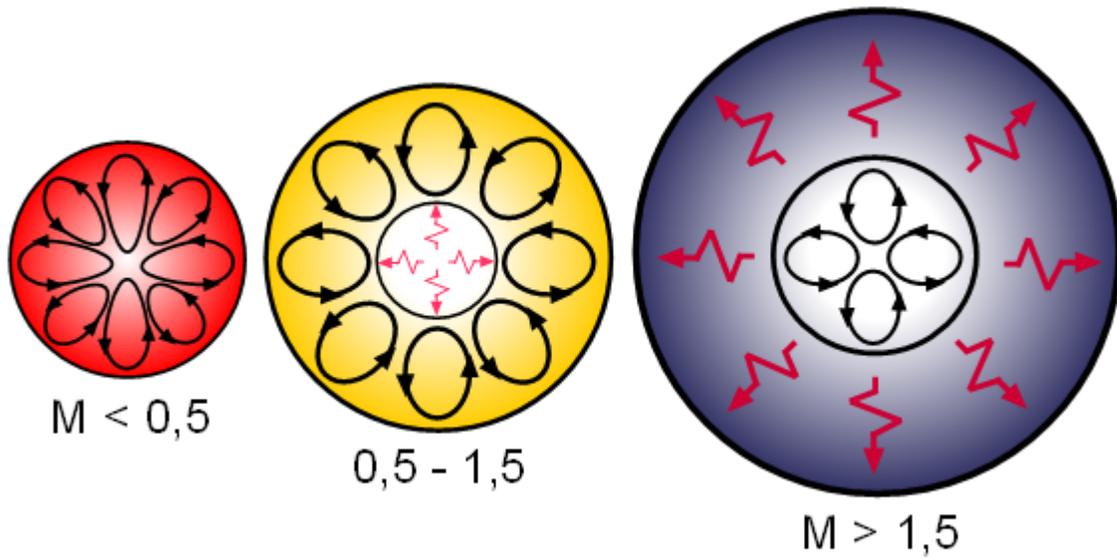Projected timeline of the Sun's life

# Birth



A dense starfield in Sagittarius.

Stellar evolution begins with the gravitational collapse of a giant molecular cloud (GMC). Typical GMCs are roughly 100 light-years ($9.5 \times 10^{14}$ km) across and contain up to 6,000,000 solar masses ($1.2 \times 10^{37}$ kg). As it collapses, a GMC breaks into smaller and smaller pieces. In each of these fragments, the collapsing gas releases gravitational potential energy as heat. As its temperature and pressure increase, a fragment condenses into a rotating sphere of superhot gas known as a protostar.

Protostars with masses less than roughly 0.08 $M_\odot$ ($1.6 \times 10^{29}$ kg) never reach temperatures high enough for nuclear fusion of hydrogen to begin. These are known as brown dwarfs. Brown dwarfs heavier than 13 Jupiter masses ($2.5 \times 10^{28}$ kg) do fuse deuterium, and some astronomers prefer to call only these objects brown dwarfs, classifying anything larger than a planet but smaller than this a sub-stellar object. Both types, deuterium-burning or not, shine dimly and die away slowly, cooling gradually over hundreds of millions of years.

For a more massive protostar, the core temperature will eventually reach 10 million kelvins, initiating the proton-proton chain reaction and allowing hydrogen to fuse, first to deuterium and then to helium. In stars of slightly over 1 $M_\odot$ ($2.0 \times 10^{30}$ kg), the CNO cycle contributes a considerable portion of the energy generation. The onset of nuclear fusion leads relatively quickly to a hydrostatic equilibrium in which energy released by the core exerts a "radiation pressure" balancing the weight of the star's matter, preventing further gravitational collapse. The star thus evolves rapidly to a stable state, beginning the main sequence phase of its evolution.

A new star will fall at a specific point on the main sequence of the Hertzsprung-Russell diagram, with the main sequence spectral type depending upon the mass of the star. Small, relatively cold, low mass red dwarfs burn hydrogen slowly and will remain on the main sequence for hundreds of billions of years, while massive, hot supergiants will leave the main sequence after just a few million years. A mid-sized star like the Sun will remain on the main sequence for about 10 billion years. The Sun is thought to be in the middle of its lifespan; thus, it is on the main sequence.

Internal structures of main sequence stars, convection zones with arrowed cycles and radiative zones with red flashes. To the left a **low-mass** red dwarf, in the center a **mid-sized** yellow dwarf and at the right a **massive** blue-white main sequence star.

# Maturity

The continuous fusion of hydrogen into helium will cause a build-up of helium in the core. The rate at which this process occurs depends on the initial mass of the star and ranges from millions to billions of years. Larger, hotter stars produce helium more rapidly than smaller, cooler ones.

The accumulation of helium in the core causes a gradual increase in the rate of fusion and gravitational self-compression, as helium is denser than hydrogen. Higher temperatures must be attained to resist this increase in gravitational compression and to maintain a steady state.

Eventually, the core exhausts its supply of hydrogen, and without the outward pressure generated by the fusion of hydrogen to counteract the force of gravity, it contracts until either electron degeneracy becomes sufficient to oppose gravity or the core becomes hot enough (around 100 megakelvins) for helium fusion to begin. Which of these happens first depends upon the star's mass.

## Low-mass stars

What happens after a low-mass star ceases to produce energy through fusion is not directly known: the universe is thought to be around 13.7 billion years old, which is less time (by several orders of magnitude, in some cases) than it takes for the fusion to cease in such stars. Current theory is based on computer modelling done by astronomers such as Don VandenBerg.

A star of less than about 0.5 solar mass will never be able to fuse helium even after the core ceases hydrogen fusion. There simply is not a stellar envelope massive enough to exert enough pressure on the core. These are the red dwarfs, such as Proxima Centauri, some of which will live thousands of times longer than the Sun. Recent astrophysical models suggest that red dwarfs of 0.1 solar mass may stay on the main sequence for almost six trillion years, and take several hundred billion more to slowly collapse into a white dwarf. If a star's core becomes stagnant (as is thought will be the case for the Sun), it will still be surrounded by layers of hydrogen which the star may subsequently draw upon. However, if the star is fully convective (as thought to be the case for the lowest-mass stars), it will not have such surrounding layers. If it does, it will develop into a red giant as described for mid-sized stars below, but never fuse helium as they do; otherwise, it will simply contract until electron degeneracy pressure halts its collapse, thus directly turning into a white dwarf.
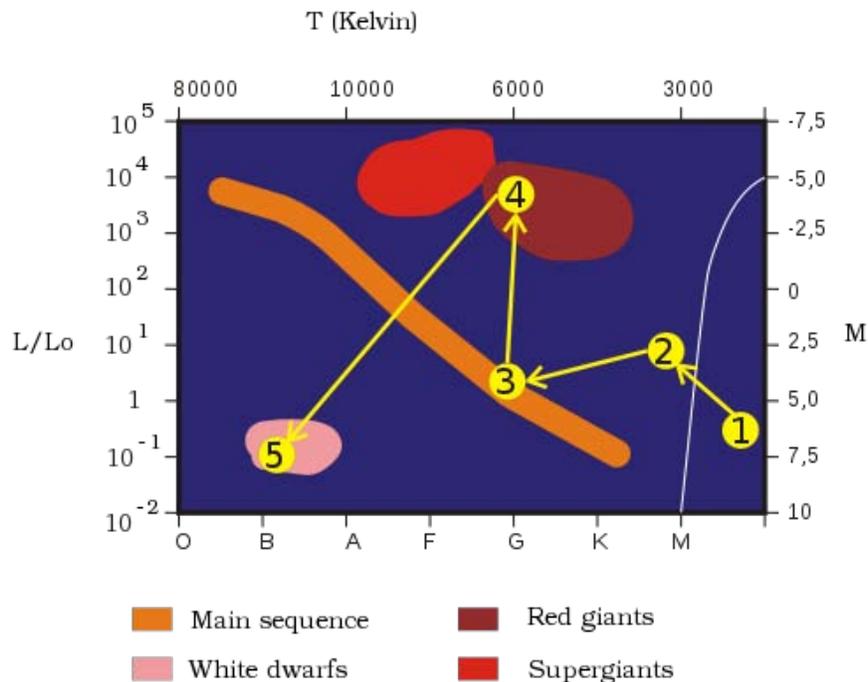
## Mid-sized stars



The Cat's Eye Nebula, a planetary nebula formed by the death of a star with about the same mass as the Sun.

In either case of the star beginning helium fusion or halting fusion due to hydrostatic equilibrium from electron degeneracy pressure, the accelerated fusion in the hydrogen-containing layer immediately over the core causes the star to expand. This lifts the outer layers away from the core, reducing the gravitational pull on them, and they expand faster than the energy production increases. This causes them to cool, which causes the star to become redder than when it was on the main sequence. Such stars are known as red giants.

According to the Hertzsprung-Russell diagram, a **red giant** is a large non-main sequence star of stellar classification K or M. Examples include Aldebaran in the constellation Taurus and Arcturus in the constellation of Boötes.

A star of up to a few solar masses will develop a helium core supported by electron degeneracy pressure, surrounded by layers which still contain hydrogen. Its gravity compresses the hydrogen in the layer immediately above it, causing it to fuse faster than hydrogen would fuse in a main-sequence star of the same mass. This in turn causes the star to become more luminous (from 1,000–10,000 times brighter) and expand; the degree of expansion outstrips the increase in luminosity, causing the effective temperature to decrease.

The expanding outer layers of the star are convective, with the material being mixed by turbulence from near the fusing regions up to the surface of the star. For all but the lowest-mass stars, the fused material has remained deep in the stellar interior prior to this point, so the convecting envelope makes fusion products visible at the star's surface for the first time. At this stage of evolution, the results are subtle, with the largest effects, alterations to the isotopes of hydrogen and helium, being unobservable. The effects of the CNO cycle appear at the surface, with lower $^{12}C/^{13}C$ ratios and altered proportions of carbon and nitrogen. These are detectable with spectroscopy and have been measured for many evolved stars.



Simplified illustration of the evolution of a star with the mass of the Sun.
The star forms from a collapsing cloud of gas (1),

and then undergoes a contraction period as a protostar (2),
before joining the main sequence (3).
Once the Hydrogen at the core is consumed it expands into a red giant (4),
then sheds its envelope into a planetary nebula and degenerates into a white dwarf (5).

As the hydrogen around the core is consumed, the core absorbs the resulting helium, causing it to contract further, which in turn causes the remaining hydrogen to fuse even faster. This eventually leads to ignition of helium fusion (which includes the triple-alpha process) in the core. In stars of more than approximately 0.5 solar masses, electron degeneracy pressure may delay helium fusion for millions or tens of millions of years; in more massive stars, the combined weight of the helium core and the overlying layers means that such pressure is not sufficient to delay the process significantly.
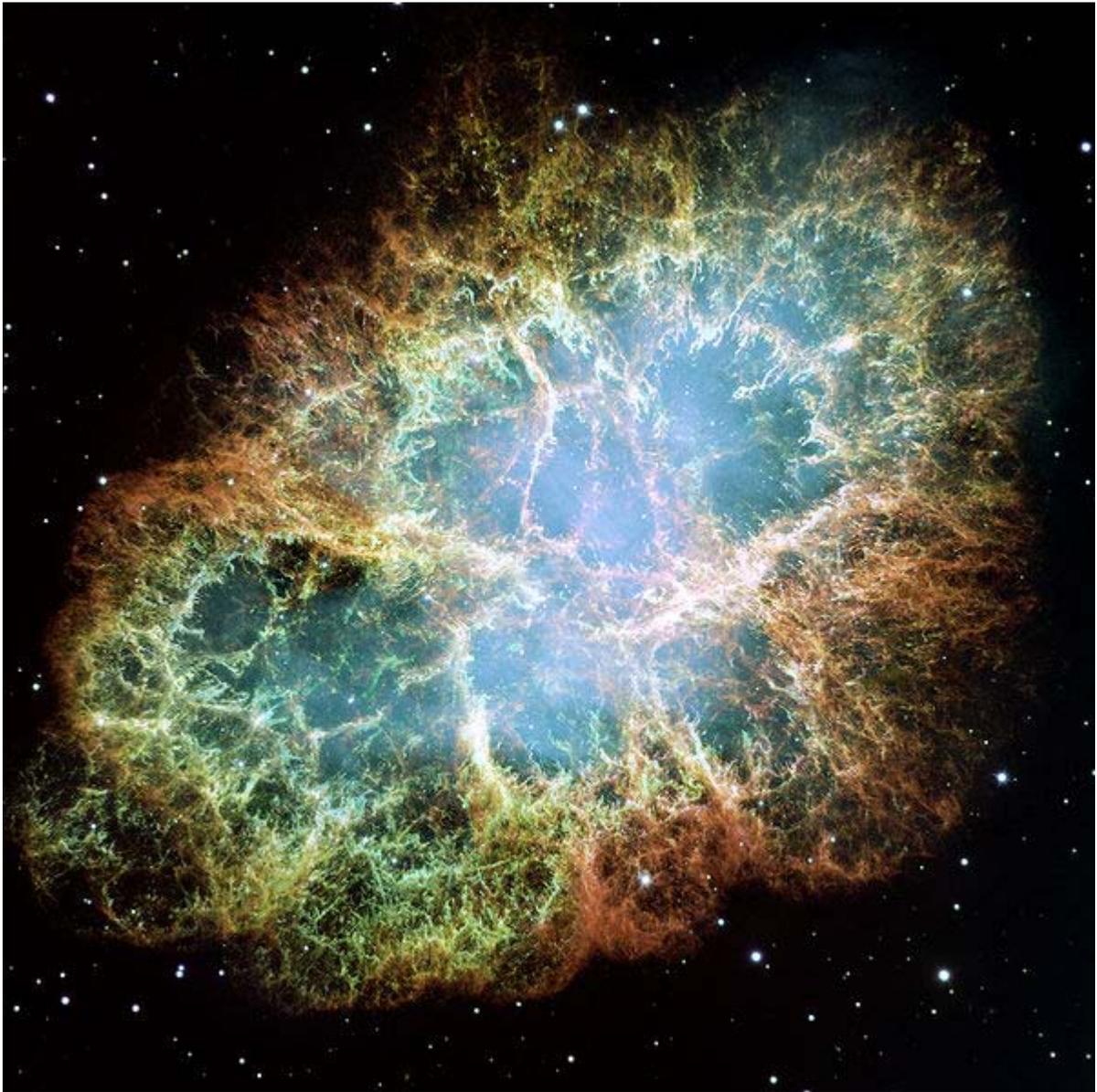
When the temperature and pressure in the core become sufficient to ignite helium fusion in the core, a helium flash will occur if the core is largely supported by electron degeneracy pressure. In more massive stars, whose core is not overwhelmingly supported by electron degeneracy pressure, the ignition of helium fusion occurs relatively quietly. Even if a helium flash does occur, the time of very rapid energy release (on the order of $10^8$ Suns) is brief, so that the visible outer layers of the star are relatively undisturbed. The energy released by helium fusion causes the core to expand, so that hydrogen fusion in the overlying layers slows and total energy generation decreases. The star contracts, although not all the way to the main sequence, and it migrates to the horizontal branch on the HR-diagram, gradually shrinking in radius and increasing its surface temperature.

After the star has consumed the helium at the core, fusion continues in a shell around a hot core of carbon and oxygen. The star follows the Asymptotic Giant Branch on the HR-diagram, paralleling the original red giant evolution, but with even faster energy generation (which lasts for a shorter time).

Changes in the energy output cause the star to change in size and temperature for certain periods. The energy output itself is shifted to lower frequency emission. This is accompanied by increased mass loss through powerful stellar winds and violent pulsations. Stars in this phase of life are called *Late type stars*, *OH-IR stars* or *Mira-type stars*, depending on their exact characteristics. The expelled gas is relatively rich in heavy elements created within the star and may be particularly oxygen or carbon enriched, depending on the type of the star. The gas builds up in an expanding shell called a circumstellar envelope and cools as it moves away from the star, allowing dust particles and molecules to form. With the high infrared energy input from the central star ideal conditions are formed in these circumstellar envelopes for maser excitation.

Helium burning reactions are extremely sensitive to temperature, which causes great instability. Huge pulsations build up and eventually give the outer layers of the star enough kinetic energy to be ejected, potentially forming a planetary nebula. At the center of the nebula remains the core of the star, which cools down to become a small but dense white dwarf.
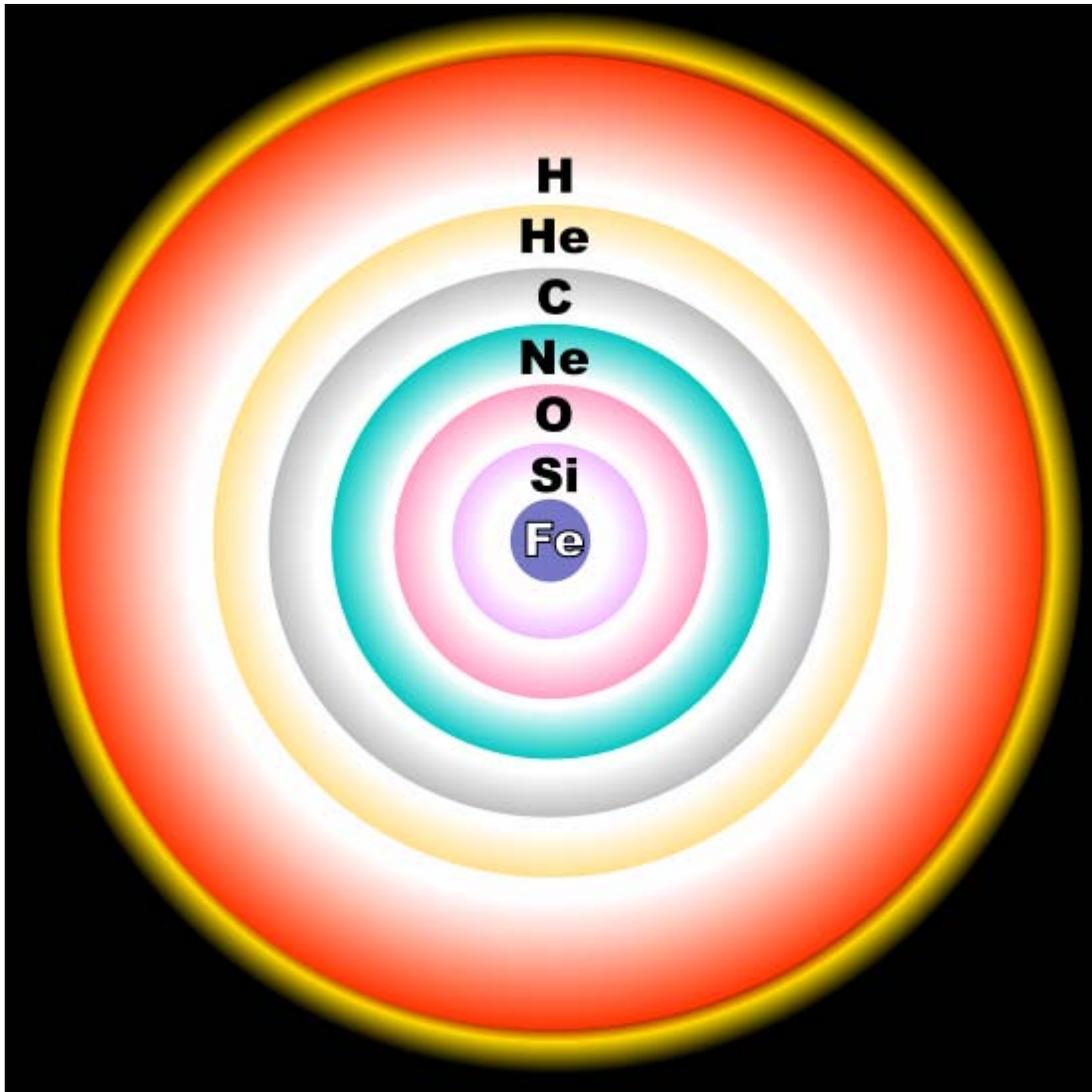
**Massive stars**



The Crab Nebula, the shattered remnants of a star which exploded as a supernova in 1054 AD.

In massive stars, the core is already large enough at the onset of hydrogen burning shell that helium ignition will occur before electron degeneracy pressure has a chance to become prevalent. Thus, when these stars expand and cool, they do not brighten as much as lower mass stars; however, they were much brighter than lower mass stars to begin with, and are thus still brighter than the red giants formed from less massive stars. These stars are unlikely to survive as red supergiants, instead they will destroy themselves as Supernovae Type II.

Extremely massive stars (more than approximately 40 solar masses), which are very luminous and thus have very rapid stellar winds, lose mass so rapidly due to radiation pressure that they tend to strip off their own envelopes before they can expand to become red supergiants, and thus retain extremely high surface temperatures (and blue-white color) from their main sequence time onwards. Stars cannot be more than about 120 solar masses because the outer layers would be expelled by the extreme radiation. Although lower mass stars normally do not burn off their outer layers so rapidly, they can likewise avoid becoming red giants or red supergiants if they are in binary systems close enough so that the companion star strips off the envelope as it expands, or if they rotate rapidly enough so that convection extends all the way from the core to the surface, resulting in the absence of a separate core and envelope due to thorough mixing.

The core grows hotter and denser as it gains material from fusion of hydrogen at the base of the envelope. In all massive stars, electron degeneracy pressure is insufficient to halt collapse by itself, so as each major element is consumed in the center, progressively heavier elements ignite, temporarily halting collapse. If the core of the star is not too massive (less than approximately 1.4 solar masses, taking into account mass loss that has occurred by this time), it may then form a white dwarf (possibly surrounded by a planetary nebula) as described above for less massive stars, with the difference that the white dwarf is composed chiefly of oxygen, neon, and magnesium.

The onion-like layers of a massive, evolved star just before core collapse. (Not to scale.)

Above a certain mass (estimated at approximately 2.5 solar masses, within the star's progenitor originally being around 10 solar masses), the core will reach the temperature (approximately 1.1 gigakelvins) at which neon partially breaks down to form oxygen and helium, the latter of which immediately fuses with some of the remaining neon to form magnesium; then oxygen fuses to form sulfur, silicon, and smaller amounts of other elements. Finally, the temperature gets high enough that any nucleus can be partially broken down, most commonly releasing an alpha particle (helium nucleus) which immediately fuses with another nucleus, so that several nuclei are effectively rearranged into a smaller number of heavier nuclei, with net release of energy because the addition of fragments to nuclei exceeds the energy required to break them off the parent nuclei.

A star with a core mass too great to form a white dwarf but insufficient to achieve sustained conversion of neon to oxygen and magnesium, will undergo core collapse (due to electron capture) before achieving fusion of the heavier elements. Both heating and

cooling caused by electron capture onto minor constituent elements (such as aluminum and sodium) prior to collapse may have a significant impact on total energy generation within the star shortly before collapse. This may produce a noticeable effect on the abundance of elements and isotopes ejected in the subsequent supernova.

Once the nucleosynthesis process arrives at iron-56, the continuation of this process consumes energy (the addition of fragments to nuclei releases less energy than required to break them off the parent nuclei). If the mass of the core exceeds the Chandrasekhar limit, electron degeneracy pressure will be unable to support its weight against the force of gravity, and the core will undergo sudden, catastrophic collapse to form a neutron star or (in the case of cores that exceed the Tolman-Oppenheimer-Volkoff limit), a black hole. Through a process that is not completely understood, some of the gravitational potential energy released by this core collapse is converted into a Type Ib, Type Ic, or Type II supernova. It is known that the core collapse produces a massive surge of neutrinos, as observed with supernova SN 1987A. The extremely energetic neutrinos fragment some nuclei; some of their energy is consumed in releasing nucleons, including neutrons, and some of their energy is transformed into heat and kinetic energy, thus augmenting the shock wave started by rebound of some of the infalling material from the collapse of the core. Electron capture in very dense parts of the infalling matter may produce additional neutrons. As some of the rebounding matter is bombarded by the neutrons, some of its nuclei capture them, creating a spectrum of heavier-than-iron material including the radioactive elements up to (and likely beyond) uranium. Although non-exploding red giant stars can produce significant quantities of elements heavier than iron using neutrons released in side reactions of earlier nuclear reactions, the abundance of elements heavier than iron (and in particular, of certain isotopes of elements that have multiple stable or long-lived isotopes) produced in such reactions is quite different from that produced in a supernova. Neither abundance alone matches that found in the Solar System, so both supernovae and ejection of elements from red giant stars are required to explain the observed abundance of heavy elements and isotopes thereof.

The energy transferred from collapse of the core to rebounding material not only generates heavy elements, but (by a mechanism which is not fully understood) provides for their acceleration well beyond escape velocity, thus causing a Type Ib, Type Ic, or Type II supernova. Note that current understanding of this energy transfer is still not satisfactory; although current computer models of Type Ib, Type Ic, and Type II supernovae account for part of the energy transfer, they are not able to account for enough energy transfer to produce the observed ejection of material. Some evidence gained from analysis of the mass and orbital parameters of binary neutron stars (which require two such supernovae) hints that the collapse of an oxygen-neon-magnesium core may produce a supernova that differs observably (in ways other than size) from a supernova produced by the collapse of an iron core.

The most massive stars may be completely destroyed by a supernova with an energy greatly exceeding its gravitational binding energy. This rare event, caused by pair-instability, leaves behind no black hole remnant.

# Stellar remnants

After a star has burned out its fuel supply, its remnants can take one of three forms, depending on the mass during its lifetime.

## White dwarfs

For a star of 1 solar mass, the resulting white dwarf is of about 0.6 solar mass, compressed into approximately the volume of the Earth. White dwarfs are stable because the inward pull of gravity is balanced by the degeneracy pressure of the star's electrons. (This is a consequence of the Pauli exclusion principle.) Electron degeneracy pressure provides a rather soft limit against further compression; therefore, for a given chemical composition, white dwarfs of higher mass have a smaller volume. With no fuel left to burn, the star radiates its remaining heat into space for billions of years.

The chemical composition of the white dwarf depends upon its mass. A star of a few solar masses will ignite carbon fusion to form magnesium, neon, and smaller amounts of other elements, resulting in a white dwarf composed chiefly of oxygen, neon, and magnesium, provided that it can lose enough mass to get below the Chandrasekhar limit (see below), and provided that the ignition of carbon is not so violent as to blow the star apart in a supernova. A star of mass on the order of magnitude of the Sun will be unable to ignite carbon fusion, and will produce a white dwarf composed chiefly of carbon and oxygen, and of mass too low to collapse unless matter is added to it later (see below). A star of less than about half the mass of the Sun will be unable to ignite helium fusion (as noted earlier), and will produce a white dwarf composed chiefly of helium.

In the end, all that remains is a cold dark mass sometimes called a black dwarf. However, the universe is not old enough for any black dwarf stars to exist yet.

If the white dwarf's mass increases above the Chandrasekhar limit, which is 1.4 solar masses for a white dwarf composed chiefly of carbon, oxygen, neon, and/or magnesium, then electron degeneracy pressure fails due to electron capture and the star collapses. Depending upon the chemical composition and pre-collapse temperature in the center, this will lead either to collapse into a neutron star or runaway ignition of carbon and oxygen. Heavier elements favor continued core collapse, because they require a higher temperature to ignite, because electron capture onto these elements and their fusion products is easier; higher core temperatures favor runaway nuclear reaction, which halts core collapse and leads to a Type Ia supernova. These supernovae may be many times brighter than the Type II supernova marking the death of a massive star, even though the latter has the greater total energy release. This inability to collapse means that no white dwarf more massive than approximately 1.4 solar masses can exist (with a possible minor exception for very rapidly spinning white dwarfs, whose centrifugal force due to rotation partially counteracts the weight of their matter). Mass transfer in a binary system may cause an initially stable white dwarf to surpass the Chandrasekhar limit.

If a white dwarf forms a close binary system with another star, hydrogen from the larger companion may accrete around and onto a white dwarf until it gets hot enough to fuse in a runaway reaction at its surface, although the white dwarf remains below the Chandrasekhar limit. Such an explosion is termed a nova.

## Neutron stars



Bubble-like shock wave still expanding from a supernova explosion 15,000 years ago.

A **neutron star** is a type of remnant that can result from the gravitational collapse of a massive star during a Type II, Type Ib or Type Ic supernova event. Such stars are composed almost entirely of neutrons, which are subatomic particles without electrical charge and a slightly larger mass than protons. Neutron stars are very hot and are supported against further collapse because of the Pauli exclusion principle. This principle states that no two neutrons (or any other fermionic particle) can occupy the same place and quantum state simultaneously.
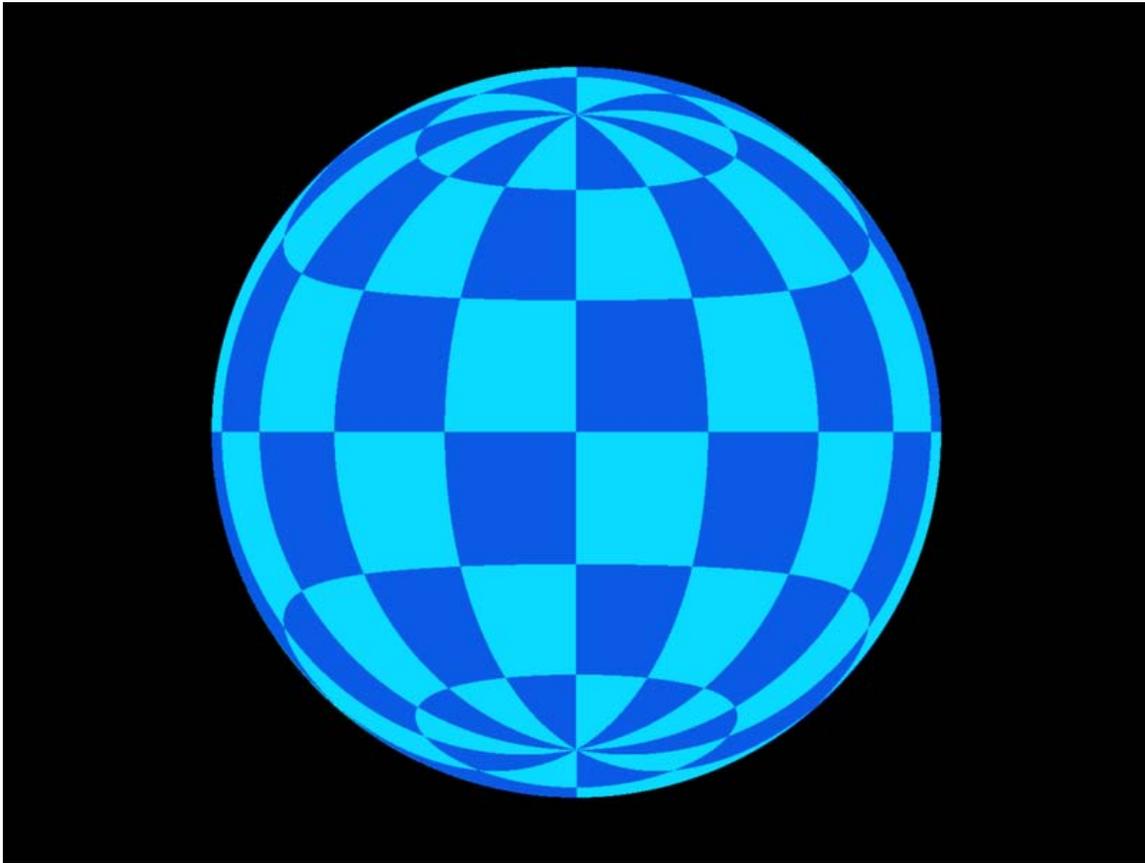
A typical neutron star has a mass between 1.35 and about 2.1 solar masses, with a corresponding radius of about 12 km if the Akmal-Pandharipande-Ravenhall (APR) Equation of state (EOS) is used. In contrast, the Sun's radius is about 60,000 times that. Neutron stars have overall densities predicted by the APR EOS of $3.7 \times 10^{17}$ to $5.9 \times 10^{17}$ kg/m$^3$ ($2.6 \times 10^{14}$ to $4.1 \times 10^{14}$ times the density of the Sun), which compares with the approximate density of an atomic nucleus of $3 \times 10^{17}$ kg/m$^3$. The neutron star's density varies from below $1 \times 10^9$ kg/m$^3$ in the crust increasing with depth to above $6 \times 10^{17}$ or $8 \times 10^{17}$ kg/m$^3$ deeper inside. This density is approximately equivalent to the mass of the entire human population compressed to the size of a sugar cube.

In general, compact stars of less than 1.44 solar masses – the Chandrasekhar limit – are white dwarfs, and above 2 to 3 solar masses (the Tolman-Oppenheimer-Volkoff limit), a quark star might be created; however, this is uncertain. Gravitational collapse will usually occur on any compact star between 10 and 25 solar masses and produce a black hole.

# Formation

As the core of a massive star is compressed during a supernova, and collapses into a neutron star, it retains most of its angular momentum. Since it has only a tiny fraction of its parent's radius (and therefore its moment of inertia is sharply reduced), a neutron star is formed with very high rotation speed, and then gradually slows down. Neutron stars are known to have rotation periods between about 1.4 ms to 30 seconds. The neutron star's compactness also gives it very high surface gravity, up to $7 \times 10^{12}$ m/s² with typical values of a few $\times 10^{12}$ m/s² (that is more than $10^{11}$ times of that of Earth). One measure of such immense gravity is the fact that neutron stars have an escape velocity of around 100,000 km/s, about 33% of the speed of light. Matter falling onto the surface of a neutron star would be accelerated to tremendous speed by the star's gravity. The force of impact would likely destroy the object's component atoms, rendering all its matter identical, in most respects, to the rest of the star.

# Properties



Gravitational light deflection at a neutron star. Due to relativistic light deflection more than half of the surface is visible (each chequered patch here represents 30 degrees by 30 degrees). The mass of the star depicted here is 1 and its radius 4, in natural units from a geometrized unit system such that it has double its Schwarzschild radius of 2.

The gravitational field at the star's surface is about $2 \times 10^{11}$ times stronger than on Earth. The escape velocity is about 100,000 km/s, which is about one third the speed of light. Such a strong gravitational field acts as a gravitational lens and bends the radiation emitted by the star such that parts of the normally invisible rear surface become visible.

The gravitational binding energy of a neutron star with two solar masses is equivalent to the total conversion of one solar mass to energy (from the law of mass-energy equivalence, $E = mc^2$). That energy was released during the supernova explosion.
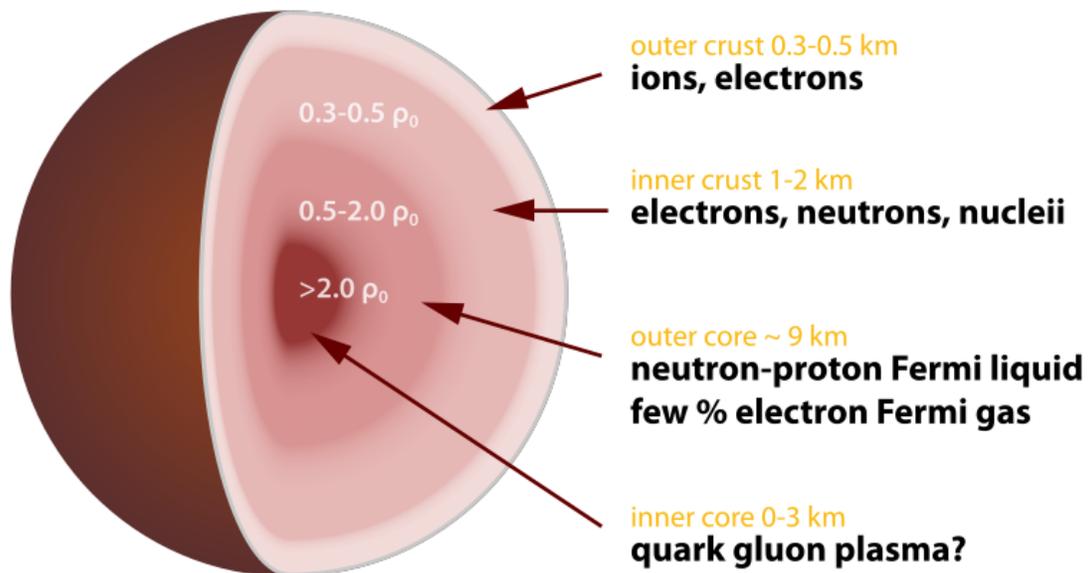
A neutron star is so dense that one teaspoon (5 milliliters) of its material would have a mass over $5.5 \times 10^{12}$ kg, about 900 times the mass of the Great Pyramid of Giza. The resulting force of gravity is so strong that if an object were to fall from a height of one

meter it would only take one microsecond to hit the surface of the neutron star, and would do so at around 2000 kilometers per second, or 7.2 million kilometers per hour.

The temperature inside a newly formed neutron star is from around $10^{11}$ to $10^{12}$ kelvin. However, the huge number of neutrinos it emits carries away so much energy that the temperature falls within a few years to around $10^6$ kelvin. Even at 1 million kelvin, most of the light generated by a neutron star is in X-rays. In visible light, neutron stars probably radiate approximately the same energy in all parts of visible spectrum, and therefore appear white.

The equation of state (EOS) for a neutron star is still not known. It is assumed that it differs significantly from that of a white dwarf, whose EOS is that of a degenerate gas which can be described in close agreement with special relativity. However, with a neutron star the increased effects of general relativity can no longer be ignored. Several EOS have been proposed (FPS, UU, APR, L, SLy, and others) and current research is still attempting to constrain the theories to make predictions of neutron star matter. This means that the relation between density and mass is not fully known, and this causes uncertainties in radius estimates. For example, a 1.5 solar mass neutron star could have a radius of 10.7, 11.1, 12.1 or 15.1 kilometres (for EOS FPS, UU, APR or L respectively). All EOS show that neutronium compresses with pressure.

# Structure



Cross-section of neutron star. Densities are in terms of $\rho_0$ the saturation nuclear matter density, where nucleons begin to touch.

Current understanding of the structure of neutron stars is defined by existing mathematical models, but it might be possible to infer through studies of neutron-star oscillations. Similar to asteroseismology for ordinary stars, the inner structure might be derived by analyzing observed frequency spectra of stellar oscillations.

On the basis of current models, the matter at the surface of a neutron star is composed of ordinary atomic nuclei crushed into a solid lattice with a sea of electrons flowing through the gaps between them. It is possible that the nuclei at the surface are iron, due to iron's high binding energy per nucleon. It is also possible that heavy element cores, such as iron, simply drown beneath the surface, leaving only light nuclei like helium and hydrogen cores. If the surface temperature exceeds $10^6$ kelvin (as in the case of a young pulsar), the surface should be fluid instead of the solid phase observed in cooler neutron stars (temperature $<10^6$ kelvin).
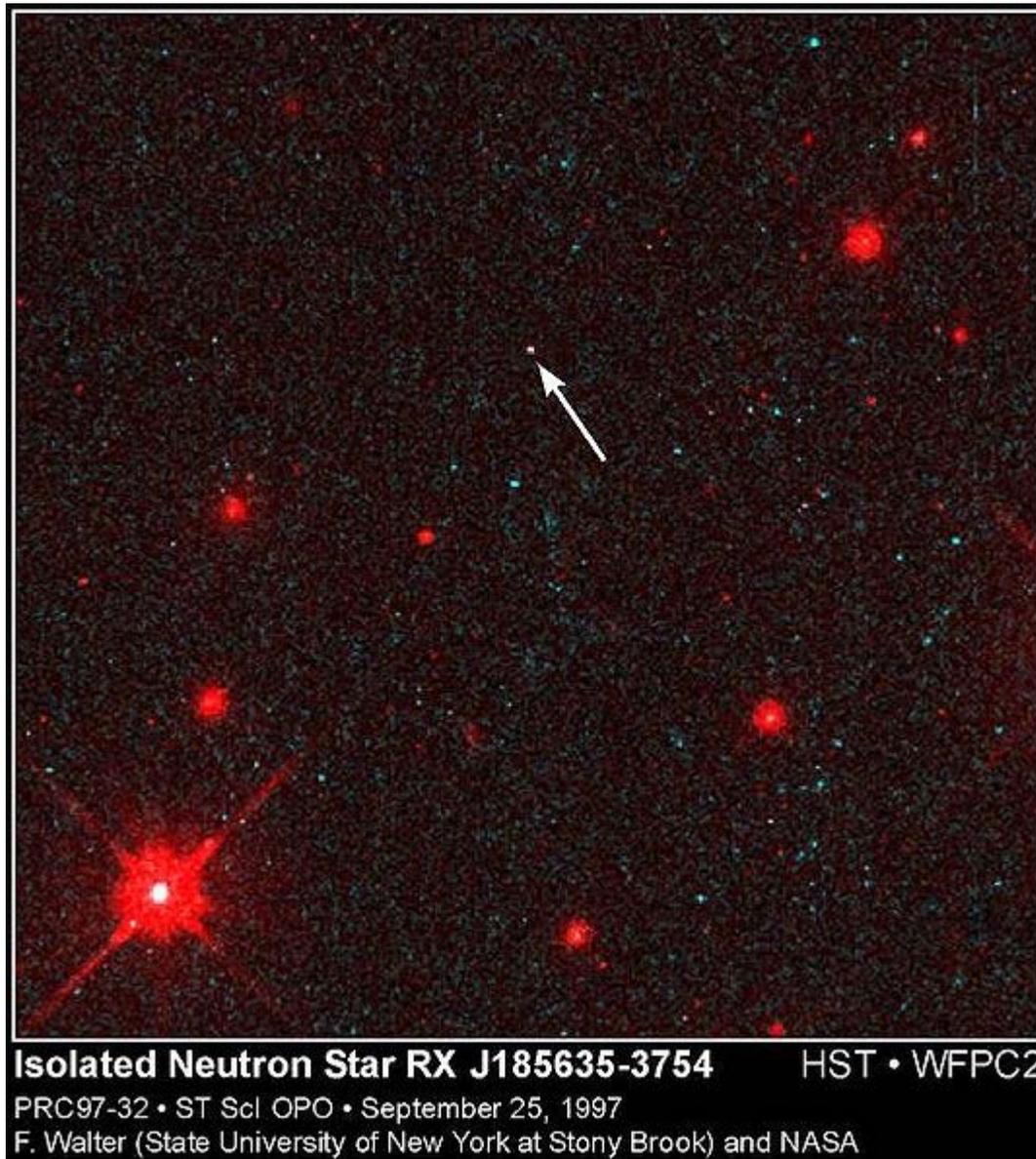
The "atmosphere" of the star is roughly one meter thick, and its dynamic is fully controlled by the star's magnetic field. Below the atmosphere one encounters a solid "crust". This crust is extremely hard and very smooth (with maximum surface irregularities of ~5 mm), because of the extreme gravitational field.

Proceeding inward, one encounters nuclei with ever increasing numbers of neutrons; such nuclei would decay quickly on Earth, but are kept stable by tremendous pressures.

Proceeding deeper, one comes to a point called neutron drip where free neutrons leak out of nuclei. In this region, there are nuclei, free electrons, and free neutrons. The nuclei become smaller and smaller until the core is reached, by definition the point where they disappear altogether. The exact nature of the superdense matter in the core is still not well understood. While this theoretical substance is referred to as neutronium in science fiction and popular literature, the term "neutronium" is rarely used in scientific publications, due to ambiguity over its meaning. The term neutron-degenerate matter is sometimes used, though not universally as the term incorporates assumptions about the nature of neutron star core material.

Neutron star core material could be a superfluid mixture of neutrons with a few protons and electrons, or it could incorporate high-energy particles like pions and kaons in addition to neutrons, or it could be composed of strange matter incorporating quarks heavier than up and down quarks, or it could be quark matter not bound into hadrons. (A compact star composed entirely of strange matter would be called a strange star.) However, so far, observations have neither indicated nor ruled out such exotic states of matter.

# History of discoveries



Isolated Neutron Star RX J185635-3754    HST • WFPC2
PRC97-32 • ST ScI OPO • September 25, 1997
F. Walter (State University of New York at Stony Brook) and NASA

The first direct observation of a neutron star in visible light. The neutron star is RX J185635-3754.

The neutron subatomic particle was discovered in 1932 by Sir James Chadwick. By bombarding the hydrogen atoms in paraffin with emissions from beryllium that was itself being bombarded with alpha particles, he demonstrated that these emissions contained a neutral particle that had about the same mass as a proton. In 1935 he was awarded the Nobel Prize in Physics for this discovery.

In 1934, Walter Baade and Fritz Zwicky proposed the existence of the neutron star, only a year after Chadwick's discovery of the neutron. In seeking an explanation for the origin of a supernova, they proposed that the neutron star is formed in a supernova. Supernovae are suddenly appearing dying stars in the sky, whose luminosity in the optical light outshine an entire galaxy for days to weeks. Baade and Zwicky correctly proposed at that time that the release of the gravitational binding energy of the neutron stars powers the supernova: "In the supernova process mass in bulk is annihilated". If the central part of a massive star before its collapse contains (for example) 3 solar masses, then a neutron star of 2 solar masses can be formed. The binding energy $E$ of such a neutron star, when expressed in mass units via the mass-energy equivalence formula $E = mc^2$, is 1 solar mass. It is ultimately this energy that powers the supernova.

In 1965, Antony Hewish and Samuel Okoye discovered "an unusual source of high radio brightness temperature in the Crab Nebula". This source turned out to be the Crab Nebula neutron star that resulted from the great supernova of 1054.

In 1967, Iosif Shklovsky examined the X-ray and optical observations of Scorpius X-1 and correctly concluded that the radiation comes from a neutron star at the stage of accretion.

In 1967, Jocelyn Bell and Antony Hewish discovered regular radio pulses from the location of the Hewish and Okoye radio source. This pulsar was later interpreted as originating from an isolated, rotating neutron star. The energy source of the pulsar is the rotational energy of the neutron star. The largest number of known neutron stars are of this type.

In 1971, Riccardo Giacconi, Herbert Gursky, Ed Kellogg, R. Levinson, E. Schreier, and H. Tananbaum discovered 4.8 second pulsations in an X-ray source in the constellation Centaurus, Cen X-3. They interpreted this as resulting from a rotating hot neutron star. The energy source is gravitational and results from a rain of gas falling onto the surface of the neutron star from a companion star or the interstellar medium.
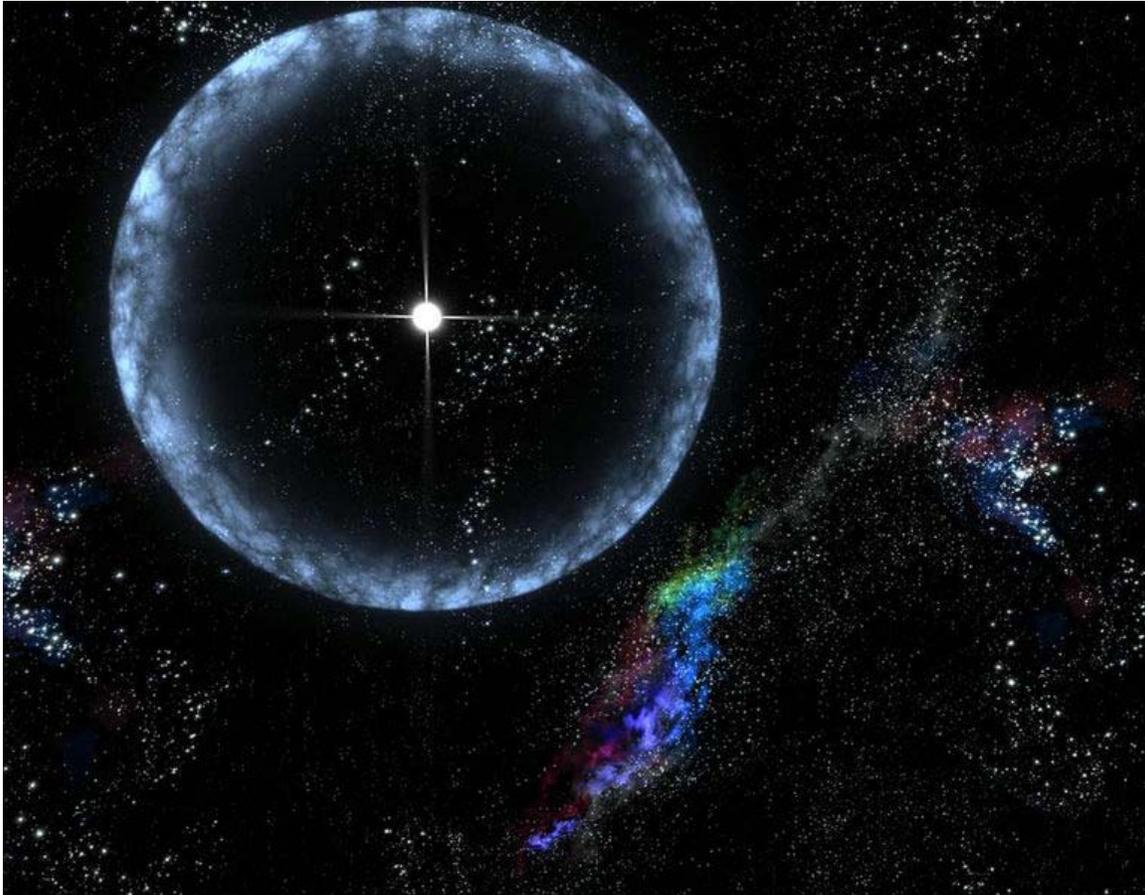
In 1974, Antony Hewish was awarded the Nobel Prize in Physics "for his decisive role in the discovery of pulsars" without Samuel Okoye and Jocelyn Bell who shared in the discovery.

# Rotation

Neutron stars rotate extremely rapidly after their creation due to the conservation of angular momentum; like spinning ice skaters pulling in their arms, the slow rotation of the original star's core speeds up as it shrinks. A newborn neutron star can rotate several times a second; sometimes, the neutron star absorbs orbiting matter from a companion star, increasing the rotation to several hundred times per second, reshaping the neutron star into an oblate spheroid.

Over time, neutron stars slow down because their rotating magnetic fields radiate energy; older neutron stars may take several seconds for each revolution.

The rate at which a neutron star slows its rotation is usually constant and very small: the observed rates of decline are between $10^{-10}$ and $10^{-21}$ seconds for each rotation. Therefore, for a typical slow down rate of $10^{-15}$ seconds per rotation, a neutron star now rotating in 1 second will rotate in 1.000003 seconds after a century, or 1.03 seconds after 1 million years.

A "starquake", or "stellar quake"

Sometimes a neutron star will *spin up* or undergo a *glitch*, a sudden small increase of its rotation speed. Glitches are thought to be the effect of a starquake - as the rotation of the star slows down, the shape becomes more spherical. Due to the stiffness of the 'neutron' crust, this happens as discrete events as the crust ruptures, similar to tectonic earthquakes. After the starquake, the star will have a smaller equatorial radius, and since angular momentum is conserved, rotational speed increases. Recent work, however, suggests that a starquake would not release sufficient energy for a neutron star glitch; it has been suggested that glitches may instead be caused by transitions of vortices in the superfluid core of the star from one metastable energy state to a lower one.

Neutron stars have been observed to "pulse" radio and x-ray emissions believed caused by particle acceleration near the magnetic poles, which need not be aligned with the rotation axis of the star. Through mechanisms not yet entirely understood, these particles produce coherent beams of radio emission. External viewers see these beams as pulses of radiation whenever the magnetic pole sweeps past the line of sight. The pulses come at the same rate as the rotation of the neutron star, and thus, appear periodic. Neutron stars which emit such pulses are called pulsars.

The most rapidly rotating neutron star currently known, PSR J1748-2446ad, rotates at 716 revolutions per second. A recent paper reported the detection of an X-ray burst oscillation (an indirect measure of spin) at 1122 Hz from the neutron star XTE J1739-285. However, at present this signal has only been seen once, and should be regarded as tentative until confirmed in another burst from this star.

## Population and distances

At present there are about 2000 known neutron stars in the Milky Way and the Magellanic Clouds, the majority of which have been detected as radio pulsars. The population of neutron stars is concentrated along the disk of the Milky Way although the spread perpendicular to the disk is fairly large. The reason for this spread is that neutron stars are born with high speeds (400 km/s) as a result of an imparted momentum-kick from an asymmetry during the supernova explosion process. One of the closest known neutron stars is PSR J0108-1431 at a distance of about 130 parsecs (or 424 light years). Another nearby neutron star that was detected transiting the backdrop of the constellation Ursa Minor has been catalogued as 1RXS J141256.0+792204. This rapidly moving object, nicknamed by its Canadian and American discoverers "Calvera", was discovered using the ROSAT/Bright Source Catalog. Initial measurements placed its distance from earth at 200 to 1,000 light years away, with later claims at about 450 light-years.

## Binary neutron stars

About 5% of all neutron stars are members of a binary system. The formation and evolution scenario of binary neutron stars is a rather exotic and complicated process. The companion stars may be either ordinary stars, white dwarfs or other neutron stars. According to modern theories of binary evolution it is expected that neutron stars also exist in binary systems with black hole companions. Such binaries are expected to be prime sources for emitting gravitational waves. Neutron stars in binary systems often emit X-rays which is caused by the heating of material (gas) accreted from the companion star. Material from the outer layers of a (bloated) companion star is sucked towards the neutron star as a result of its very strong gravitational field. As a result of this process binary neutron stars may also coalesce into black holes if the accretion of mass takes place under extreme conditions.

# Subtypes

- Neutron star
    - Protoneutron star (PNS), theorized.
    - Radio-quiet neutron stars
    - Radio loud neutron star
        - Single pulsars–general term for neutron stars that emit directed pulses of radiation towards us at regular intervals (due to their strong magnetic fields).
            - Rotation-powered pulsar *("radio pulsar")*
                - Magnetar–a neutron star with an extremely strong magnetic field (1000 times more than a regular neutron star), and long rotation periods (5 to 12 seconds).
                    - Soft gamma repeater (SGR)
                    - Anomalous X-ray pulsar (AXP)
        - Binary pulsars
            - Low-mass X-ray binaries (LMXB)
            - Intermediate-mass X-ray binaries (IMXB)
            - High-mass X-ray binaries (HMXB)
            - Accretion-powered pulsar *("X-ray pulsar")*
                - X-ray burster–a neutron star with a low mass binary companion from which matter is accreted resulting in irregular bursts of energy from the surface of the neutron star.
                - Millisecond pulsar (MSP) *("recycled pulsar")*
                    - Sub-millisecond pulsar
    - Exotic star
        - Quark star–currently a hypothetical type of neutron star composed of quark matter, or strange matter. As of 2008, there are three candidates.
        - Preon star–currently a hypothetical type of neutron star composed of preon matter. As of 2008, there is no evidence for the existence of preons.
        - Q star–currently a hypothetical type of heavy neutron star with an exotic state of matter. As of 2008, there is no evidence for their existence.
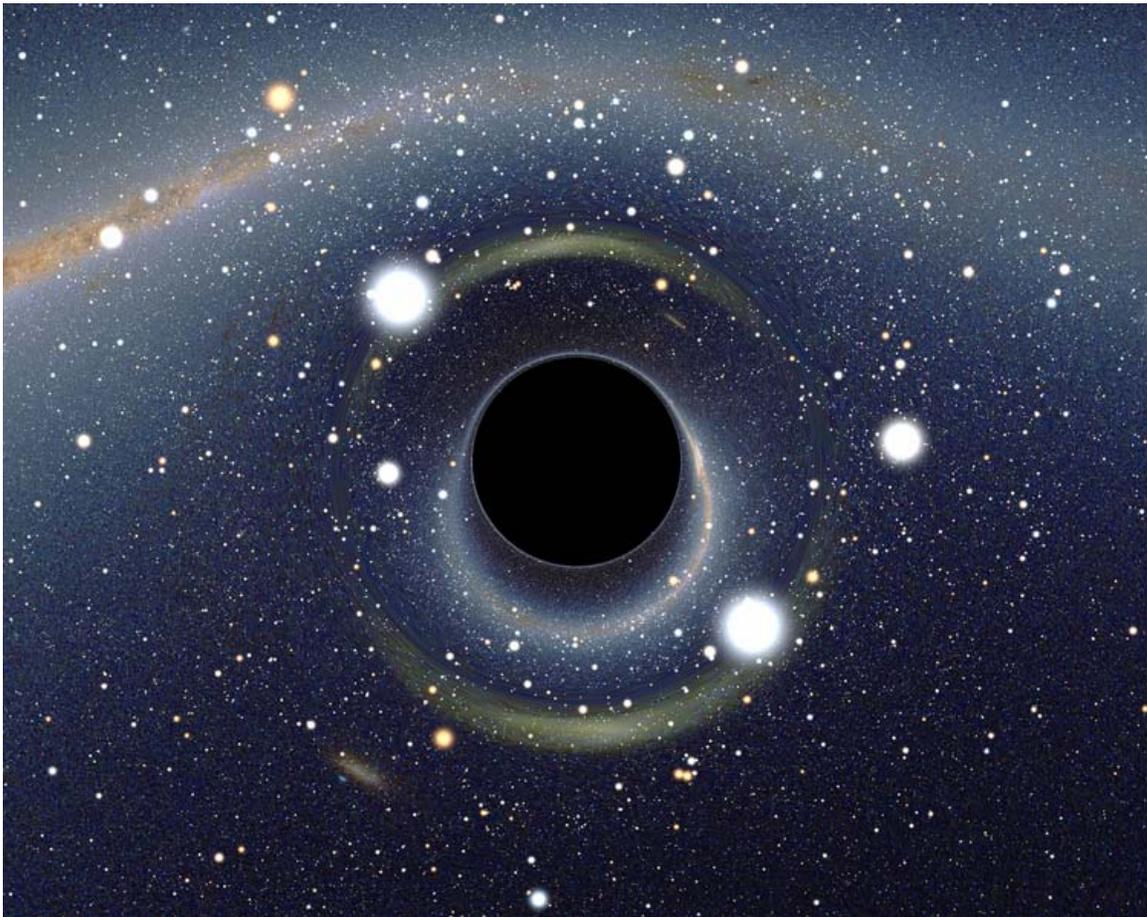
# Giant nuclei

A neutron star has some of the properties of an atomic nucleus, including density, and being made of nucleons. In popular scientific writing, neutron stars are therefore sometimes described as giant nuclei. However, in other respects, neutron stars and atomic nuclei are quite different. In particular, a nucleus is held together by the strong force,

while a neutron star is held together by gravity. It is generally more useful to consider such objects as stars.

# Examples of neutron stars

- PSR J0108-1431 - closest neutron star
- LGM-1 - the first recognized radio-pulsar
- PSR B1257+12 - the first neutron star discovered with planets (a millisecond pulsar)

**Black holes**



Simulated view of a black hole in front of the Large Magellanic Cloud. The ratio between the black hole Schwarzschild radius and the observer distance to it is 1:9. Of note is the gravitational lensing effect known as an Einstein ring, which produces a set of two fairly bright and large but highly distorted images of the Cloud as compared to its actual angular size.

A **black hole** is a region of space from which nothing, not even light, can escape. It is the result of the deformation of spacetime caused by a very compact mass. Around a black

hole there is an undetectable surface which marks the point of no return, called an event horizon. It is called "black" because it absorbs all the light that hits it, reflecting nothing, just like a perfect black body in thermodynamics. Quantum mechanics predicts that black holes also emit radiation like a black body with a finite temperature. This temperature decreases with the mass of the black hole, making it unlikely to observe this radiation for black holes of stellar mass.

Despite its invisible interior, a black hole can be observed through its interaction with other matter. A black hole can be inferred by tracking the movement of a group of stars that orbit a region in space. Alternatively, when gas falls into a stellar black hole from a companion star, the gas spirals inward, heating to very high temperatures and emitting large amounts of radiation that can be detected from earthbound and Earth-orbiting telescopes.

Astronomers have identified numerous stellar black hole candidates, and have also found evidence of supermassive black holes at the center of galaxies. In 1998, astronomers found compelling evidence that a supermassive black hole of more than 2 million solar masses is located near the Sagittarius A* region in the center of the Milky Way galaxy, and more recent results using additional data find evidence that the supermassive black hole is more than 4 million solar masses.

# History



Simulation of gravitational lensing by a black hole which distorts the image of a galaxy in the background

The idea of a body so massive that even light could not escape was first put forward by geologist John Michell in a letter written to Henry Cavendish in 1783 to the Royal Society:

If the semi-diameter of a sphere of the same density as the Sun were to exceed that of the Sun in the proportion of 500 to 1, a body falling from an infinite height towards it would have acquired at its surface greater velocity than that of light, and consequently supposing light to be attracted by the same force in proportion to its vis inertiae, with

other bodies, all light emitted from such a body would be made to return towards it by its own proper gravity.
—John Michell

In 1796, mathematician Pierre-Simon Laplace promoted the same idea in the first and second editions of his book *Exposition du système du Monde* (it was removed from later editions). Such "dark stars" were largely ignored in the nineteenth century, since it was then thought that a massless wave such as light could not be influenced by gravity.

## General relativity

In 1915, Albert Einstein developed his theory of general relativity, having earlier shown that gravity does influence light's motion. A few months later, Karl Schwarzschild gave the solution for the gravitational field of a point mass and a spherical mass. A few months after Schwarzschild, Johannes Droste, a student of Hendrik Lorentz, independently gave the same solution for the point mass and wrote more extensively about its properties. This solution had a peculiar behaviour at what is now called the Schwarzschild radius, where it became singular, meaning that some of the terms in the Einstein equations became infinite. The nature of this surface was not quite understood at the time. In 1924, Arthur Eddington showed that the singularity disappeared after a change of coordinates, although it took until 1933 for Georges Lemaître to realize that this meant the singularity at the Schwarzschild radius was an unphysical coordinate singularity.

In 1931, Subrahmanyan Chandrasekhar calculated, using general relativity, that a non-rotating body of electron-degenerate matter above 1.44 solar masses (the Chandrasekhar limit) would collapse. His arguments were opposed by many of his contemporaries like Eddington and Lev Landau, who argued that some yet unknown mechanism would stop the collapse. They were partly correct: a white dwarf slightly more massive than the Chandrasekhar limit will collapse into a neutron star, which is itself stable because of the Pauli exclusion principle. But in 1939, Robert Oppenheimer and others predicted that neutron stars above approximately three solar masses (the Tolman–Oppenheimer–Volkoff limit) would collapse into black holes for the reasons presented by Chandrasekhar, and concluded that no law of physics was likely to intervene and stop at least some stars from collapsing to black holes.

Oppenheimer and his co-authors interpreted the singularity at the boundary of the Schwarzschild radius as indicating that this was the boundary of a bubble in which time stopped. This is a valid point of view for external observers, but not for infalling observers. Because of this property, the collapsed stars were called "frozen stars," because an outside observer would see the surface of the star frozen in time at the instant where its collapse takes it inside the Schwarzschild radius. This is a known property of modern black holes, but it must be emphasized that the light from the surface of the frozen star becomes redshifted very fast, turning the black hole black very quickly. Many physicists could not accept the idea of time standing still at the Schwarzschild radius, and there was little interest in the subject for over 20 years.

## Golden age

In 1958, David Finkelstein identified the Schwarzschild surface $r = 2m$ [in geometrized units, i.e. $2Gm/c^2$] as an event horizon, "a perfect unidirectional membrane: causal influences can cross it in only one direction". This did not strictly contradict Oppenheimer's results, but extended them to include the point of view of infalling observers. Finkelstein's solution extended the Schwarzschild solution for the future of observers falling into the black hole. A complete extension had already been found by Martin Kruskal, who was urged to publish it.

These results came at the beginning of the golden age of general relativity, which is marked by general relativity and black holes becoming mainstream subjects of research. This process was helped by the discovery of pulsars in 1967, which were within a few years shown to be rapidly rotating neutron stars. Until that time, neutron stars, like black holes, were regarded as just theoretical curiosities; but the discovery of pulsars showed their physical relevance and spurred a further interest in all types of compact objects that might be formed by gravitational collapse.

In this period more general black hole solutions where found. In 1963, Roy Kerr found the exact solution for a rotating black hole. Two years later Ezra T. Newman found the axisymmetric solution for a black hole which is both rotating and electrically charged. Through the work of Werner Israel, Brandon Carter, and D. C. Robinson the no-hair theorem emerged, stating that a stationary black hole solution is completely described by the three parameters of the Kerr–Newman metric; mass, angular momentum, and electric charge.

For a long time, it was suspected that the strange features of the black hole solutions were pathological artefacts from the symmetry conditions imposed, and that the singularities would not appear in generic situations. This view was held in particular by Belinsky, Khalatnikov, and Lifshitz, who tried to prove that no singularities appear in generic solutions. However, in the late sixties Roger Penrose and Stephen Hawking used global techniques to prove that singularities are generic.

Work by James Bardeen, Jacob Bekenstein, Carter, and Hawking in the early 1970s led to the formulation of the laws of black hole mechanics. These laws describe the behaviour of a black hole in close analogy to the laws of thermodynamics by relating mass to energy, area to entropy, and surface gravity to temperature. The analogy was completed when Hawking, in 1974, showed that quantum field theory predicts that black holes should radiate like a black body with a temperature proportional to the surface gravity of the black hole.

The term "black hole" was first publicly used by John Wheeler during a lecture in 1967. Although he is usually credited with coining the phrase, he always insisted that it was suggested to him by somebody else. The first recorded use of the term is in a 1964 letter by Anne Ewing to the American Association for the Advancement of Science. After Wheeler's use of the term, it was quickly adopted in general use.

# Properties and structure

The no-hair theorem states that, once it achieves a stable condition after formation, a black hole has only three independent physical properties: mass, charge, and angular momentum. Any two black holes that share the same values for these properties, or parameters, are indistinguishable according to classical (i.e. non-quantum) mechanics.

These properties are special because they are visible from outside the black hole. For example, a charged black hole repels other like charges just like any other charged object. Similarly, the total mass inside a sphere containing a black hole can be found by using the gravitational analog of Gauss's law, the ADM mass, far away from the black hole. Likewise, the angular momentum can be measured from far away using frame dragging by the gravitomagnetic field.

When an object falls into a black hole, any information about the shape of the object or distribution of charge on it is evenly distributed along the horizon of the black hole, and is lost to outside observers. The behavior of the horizon in this situation is closely analogous to that of a conductive stretchy membrane with friction and electrical resistance, a dissipative system. This is different from other field theories like electromagnetism, which does not have any friction or resistivity at the microscopic level, because they are time-reversible. Because the black hole eventually achieves a stable state with only three parameters, there is no way to avoid losing information about the initial conditions: the gravitational and electric fields of the black hole give very little information about what went in. The information that is lost includes every quantity that cannot be measured far away from the black hole horizon, including the total baryon number, lepton number, and all the other nearly conserved pseudo-charges of particle physics. This behavior is so puzzling that it has been called the black hole information loss paradox.

## Physical properties

The simplest black hole has mass but neither electric charge nor angular momentum. These black holes are often referred to as Schwarzschild black holes after the physicist Karl Schwarzschild who discovered this solution in 1915. According to Birkhoff's theorem, it is the only vacuum solution that is spherically symmetric. This means that there is no observable difference between the gravitational field of such a black hole and that of any other spherical object of the same mass. The popular notion of a black hole "sucking in everything" in its surroundings is therefore only correct near the black hole horizon; far away, the external gravitational field is identical to that of any other body of the same mass.

Solutions describing more general black holes also exist. Charged black holes are described by the Reissner-Nordström metric, while the Kerr metric describes a rotating black hole. The most general stationary black hole solution known is the Kerr-Newman metric, which describes a black hole with both charge and angular momentum.

While the mass of a black hole can take any positive value, the charge and angular momentum are constrained by the mass. In Planck units, the total electric charge $Q$ and the total angular momentum $J$ are expected to satisfy

$$Q^2 + \left(\frac{J}{M}\right)^2 \leq M^2$$

for a black hole of mass $M$. Black holes saturating this inequality are called extremal. Solutions of Einstein's equations violating the inequality exist, but do not have a horizon. These solutions have naked singularities and are deemed *unphysical*. The cosmic censorship hypothesis rules out the formation of such singularities through the gravitational collapse of realistic matter. This is supported by numerical simulations.

Due to the relatively large strength of the electromagnetic force, black holes forming from the collapse of stars are expected to retain the nearly neutral charge of the star. Rotation, however, is expected to be a common feature of compact objects, and the black-hole candidate binary X-ray source GRS 1915+105 appears to have an angular momentum near the maximum allowed value.

| Class | Mass | Size |
|---|---|---|
| Supermassive black hole | $\sim 10^5$–$10^9$ $M_{Sun}$ | $\sim$0.001–10 AU |
| Intermediate-mass black hole | $\sim 10^3$ $M_{Sun}$ | $\sim 10^3$ km $= R_{Earth}$ |
| Stellar black hole | $\sim 10$ $M_{Sun}$ | $\sim$30 km |
| Micro black hole | up to $\sim M_{Moon}$ | up to $\sim$0.1 mm |

Black holes are commonly classified according to their mass, independent of angular momentum $J$ or electric charge $Q$. The size of a black hole, as determined by the radius of the event horizon, or Schwarzschild radius, is roughly proportional to the mass $M$ through
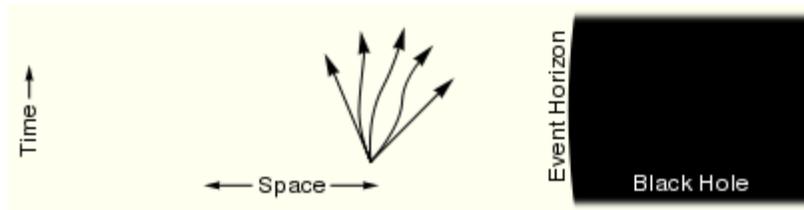
$$r_{sh} = \frac{2GM}{c^2} \approx 2.95 \, \frac{M}{M_{Sun}} \text{ km,}$$

where $r_{sh}$ is the Schwarzschild radius and $M_{Sun}$ is the mass of the Sun. This relation is exact only for black holes with zero charge and angular momentum, for more general black holes it can differ up to a factor of 2. The table on the right lists the various classes of black hole that are distinguished.

**Event horizon**



Far away from the black hole a particle can move in any direction. It is only restricted by the speed of light.



Closer to the black hole spacetime starts to deform. There are more paths going towards the black hole than paths moving away.



Inside of the event horizon all paths bring the particle closer to the center of the black hole. It is no longer possible for the particle to escape.

The defining feature of a black hole is the appearance of an event horizon—a boundary in spacetime through which matter and light can only pass inward towards the mass of the black hole. Nothing, not even light, can escape from inside the event horizon. The event horizon is referred to as such because if an event occurs within the boundary, information from that event cannot reach an outside observer, making it impossible to determine if such an event occurred.

As predicted by general relativity, the presence of a large mass deforms spacetime in such a way that the paths taken by particles bend towards the mass. At the event horizon of a black hole, this deformation becomes so strong that there are no paths that lead away from the black hole.

To a distant observer, clocks near a black hole appear to tick more slowly than those further away from the black hole. Due to this effect, known as gravitational time dilation, an object falling into a black hole appears to slow down as it approaches the event

horizon, taking an infinite time to reach it. At the same time, all processes on this object slow down causing emitted light to appear redder and dimmer, an effect known as gravitational redshift. Eventually, at a point just before it reaches the event horizon, the falling object becomes so dim that it can no longer be seen.

On the other hand, an observer falling into a black hole does not notice any of these effects as he crosses the event horizon. According to his own clock, he crosses the event horizon after a finite time, although he is unable to determine exactly when he crosses it, as it is impossible to determine the location of the event horizon from local observations.

For a non rotating (static) black hole, the Schwarzschild radius delimits a spherical event horizon. The Schwarzschild radius of an object is proportional to the mass. Rotating black holes have distorted, nonspherical event horizons. Since the event horizon is not a material surface but rather merely a mathematically defined demarcation boundary, nothing prevents matter or radiation from entering a black hole, only from exiting one. The description of black holes given by general relativity is known to be an approximation, and some scientists expect that quantum gravity effects will become significant near the vicinity of the event horizon. This would allow observations of matter near a black hole's event horizon to be used to indirectly study general relativity and proposed extensions to it.

## Singularity

At the center of a black hole as described by general relativity lies a gravitational singularity, a region where the spacetime curvature becomes infinite. For a non-rotating black hole this region takes the shape of a single point and for a rotating black hole it is smeared out to form a ring singularity lying in the plane of rotation. In both cases the singular region has zero volume. It can also be shown that the singular region contains all the mass of the black hole solution. The singular region can thus be thought of as having infinite density.

An observer falling into a Schwarzschild black hole (i.e. non-rotating and no charges) cannot avoid the singularity. Any attempt to do so will only shorten the time taken to get there. When they reach the singularity, they are crushed to infinite density and their mass is added to the total of the black hole. Before that happens, they will have been torn apart by the growing tidal forces in a process sometimes referred to as spaghettification or the noodle effect.

In the case of a charged (Reissner–Nordström) or rotating (Kerr) black hole it is possible to avoid the singularity. Extending these solutions as far as possible reveals the hypothetical possibility of exiting the black hole into a different spacetime with the black hole acting as a worm hole. The possibility of travelling to another universe is however only theoretical, since any perturbation will destroy this possibility. It also appears to be possible to follow closed timelike curves (going back to one's own past) around the Kerr singularity, which lead to problems with causality like the grandfather paradox. It is

expected that none of these peculiar effects would survive in a proper quantum mechanical treatment of rotating and charged black holes.

The appearance of singularities in general relativity is commonly perceived as signaling the breakdown of the theory. This breakdown, however, is expected; it occurs in a situation where quantum mechanical effects should describe these actions due to the extremely high density and therefore particle interactions. To date it has not been possible to combine quantum and gravitational effects into a single theory. It is generally expected that a theory of quantum gravity will feature black holes without singularities.
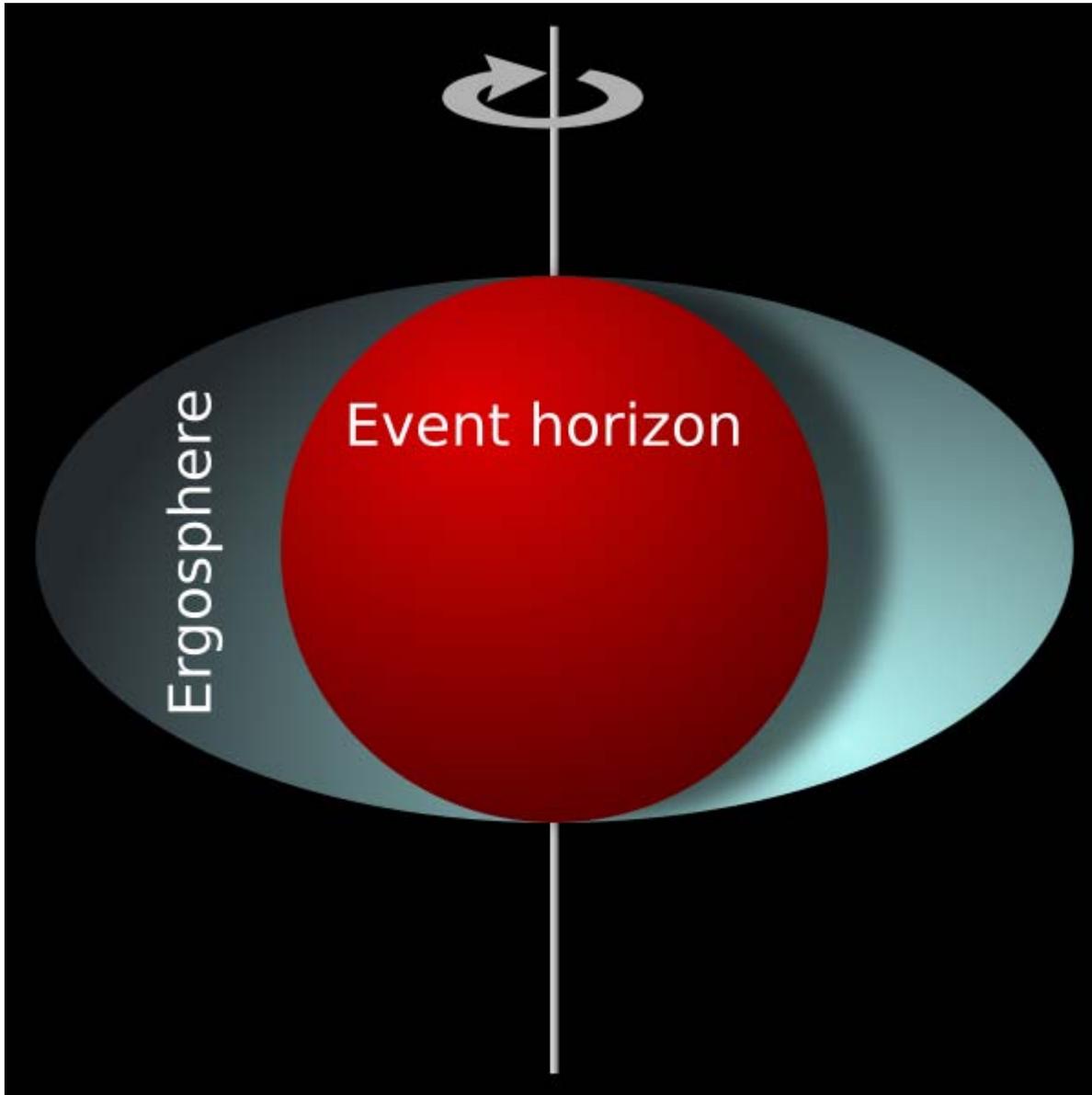
## Photon sphere

The photon sphere is a spherical boundary of zero thickness such that photons moving along tangents to the sphere will be trapped in a circular orbit. For non-rotating black holes, the photon sphere has a radius 1.5 times the Schwarzschild radius. The orbits are dynamically unstable, hence any small perturbation (such as a particle of infalling matter) will grow over time, either setting it on an outward trajectory escaping the black hole or on an inward spiral eventually crossing the event horizon.

While light can still escape from inside the photon sphere, any light that crosses the photon sphere on an inbound trajectory will be captured by the black hole. Hence any light reaching an outside observer from inside the photon sphere must have been emitted by objects inside the photon sphere but still outside of the event horizon.

Other compact objects, such as neutron stars, can also have photon spheres. This follows from the fact that the gravitational field of an object does not depend on its actual size, hence any object that is smaller than 1.5 times the Schwarzschild radius corresponding to its mass will indeed have a photon sphere.

**Ergosphere**



The ergosphere is an oblate spheroid region outside of the event horizon, where objects cannot remain stationary.

Rotating black holes are surrounded by a region of spacetime in which it is impossible to stand still, called the ergosphere. This is the result of a process known as frame-dragging; general relativity predicts that any rotating mass will tend to slightly "drag" along the spacetime immediately surrounding it. Any object near the rotating mass will tend to start moving in the direction of rotation. For a rotating black hole this effect becomes so strong

near the event horizon that an object would have to move faster than the speed of light in the opposite direction to just stand still.

The ergosphere of a black hole is bounded by the (outer) event horizon on the inside and an oblate spheroid, which coincides with the event horizon at the poles and is noticeably wider around the equator. The outer boundary is sometimes called the *ergosurface*.

Objects and radiation can escape normally from the ergosphere. Through the Penrose process, objects can emerge from the ergosphere with more energy than they entered. This energy is taken from the rotational energy of the black hole causing it to slow down.

# Formation and evolution

Considering the exotic nature of black holes, it may be natural to question if such bizarre objects could exist in nature or to suggest that they are merely pathological solutions to Einstein's equations. Einstein himself wrongly thought that black holes would not form, because he held that the angular momentum of collapsing particles would stabilize their motion at some radius. This led the general relativity community to dismiss all results to the contrary for many years. However, a minority of relativists continued to contend that black holes were physical objects, and by the end of the 1960s, they had persuaded the majority of researchers in the field that there is no obstacle to forming an event horizon.

Once an event horizon forms, Roger Penrose proved that a singularity will form somewhere inside it. Shortly afterwards, Stephen Hawking showed that many cosmological solutions describing the Big Bang have singularities without scalar fields or other exotic matter. The Kerr solution, the no-hair theorem and the laws of black hole thermodynamics showed that the physical properties of black holes were simple and comprehensible, making them respectable subjects for research. The primary formation process for black holes is expected to be the gravitational collapse of heavy objects such as stars, but there are also more exotic processes that can lead to the production of black holes.

## Gravitational collapse

Gravitational collapse occurs when an object's internal pressure is insufficient to resist the object's own gravity. For stars this usually occurs either because a star has too little "fuel" left to maintain its temperature through stellar nucleosynthesis, or because a star which would have been stable receives extra matter in a way which does not raise its core temperature. In either case the star's temperature is no longer high enough to prevent it from collapsing under its own weight (the ideal gas law explains the connection between pressure, temperature, and volume).

The collapse may be stopped by the degeneracy pressure of the star's constituents, condensing the matter in an exotic denser state. The result is one of the various types of compact star. Which type of compact star is formed depends on the mass of the remnant — the matter left over after changes triggered by the collapse (such as supernova

or pulsations leading to a planetary nebula) have blown away the outer layers. Note that this can be substantially less than the original star — remnants exceeding 5 solar masses are produced by stars which were over 20 solar masses before the collapse.

If the mass of the remnant exceeds about 3–4 solar masses (the Tolman–Oppenheimer–Volkoff limit)—either because the original star was very heavy or because the remnant collected additional mass through accretion of matter—even the degeneracy pressure of neutrons is insufficient to stop the collapse. After this, no known mechanism (except possibly quark degeneracy pressure) is powerful enough to stop the collapse and the object will inevitably collapse to a black hole.
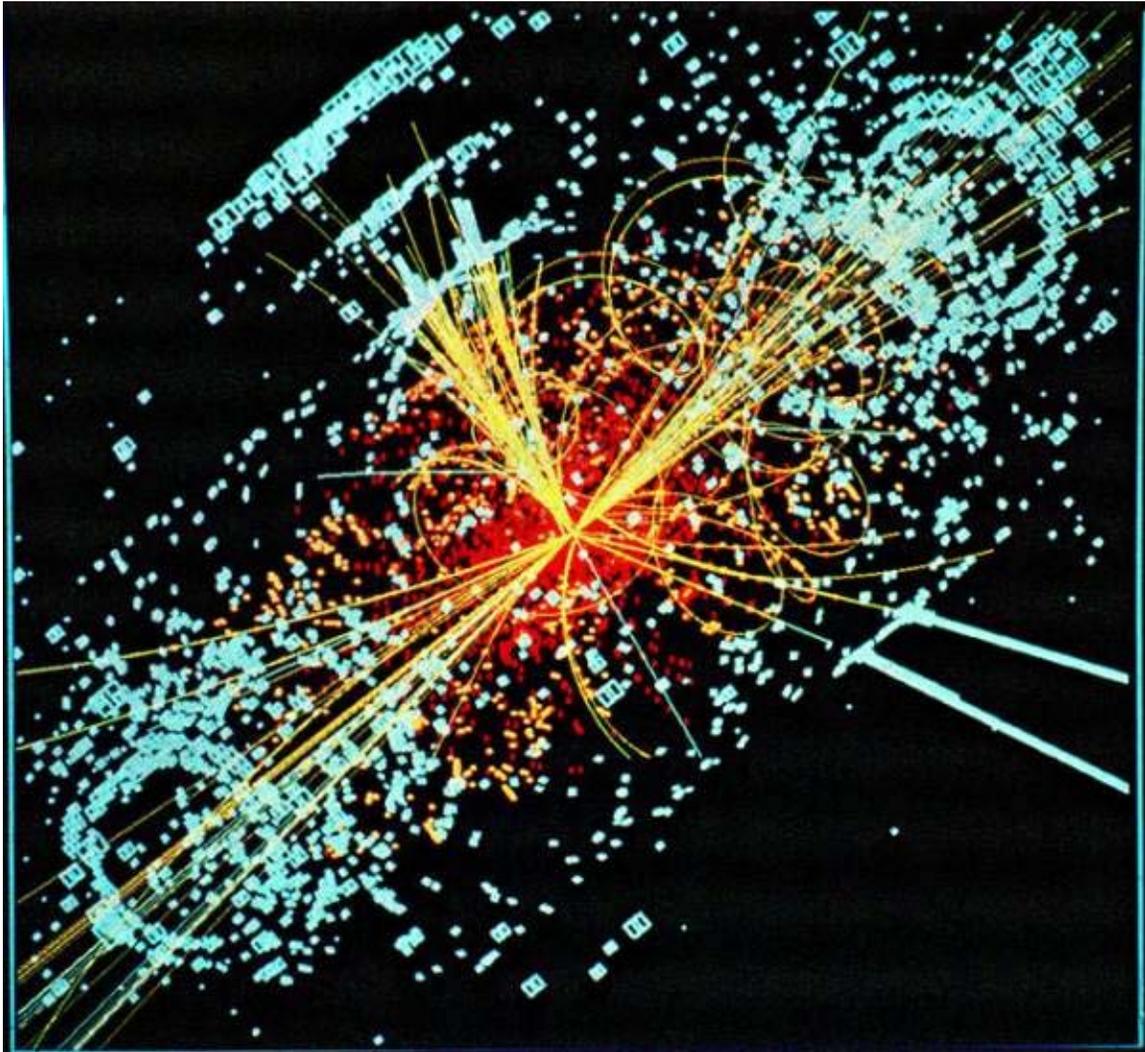
This gravitational collapse of heavy stars is assumed to be responsible for the formation of stellar mass black holes. Star formation in the young universe may have resulted in very heavy stars, which upon their collapse would have produced black holes of up to $10^3$ solar masses. These heavy black holes could be the seeds of the supermassive black holes found in the centers of most galaxies.

While most of the energy released during gravitational collapse is emitted very quickly, an outside observer does not actually see the end of this process. Even though the collapse takes a finite amount of time from the reference frame of infalling matter, a distant observer sees the infalling material slow and halt just above the event horizon, due to gravitational time dilation. Light from the collapsing material takes longer and longer to reach the observer, with the light emitted just before the event horizon forms delayed an infinite amount of time. Thus the external observer never sees the formation of the event horizon; instead, the collapsing material seems to become dimmer and increasingly red-shifted, eventually fading away.

**Primordial black holes in the Big Bang**

Gravitational collapse requires great densities. In the current epoch of the universe these high densities are only found in stars, but in the early universe shortly after the big bang densities were much greater, possibly allowing for the creation of black holes. The high density alone is not enough to allow the formation of black holes since a uniform mass distribution will not allow the mass to bunch up. In order for primordial black holes to form in such a dense medium, there must be initial density perturbations which can then grow under their own gravity. Different models for the early universe vary widely in their predictions of the size of these perturbations. Various models predict the creation of black holes, ranging from a Planck mass to hundreds of thousands of solar masses. Primordial black holes could thus account for the creation of any type of black hole.

# High-energy collisions



A simulated event in the CMS detector, a collision in which a micro black hole may be created.

Gravitational collapse is not the only process that could create black holes. In principle, black holes could also be created in high-energy collisions that create sufficient density. However, to date, no such events have ever been detected either directly or indirectly as a deficiency of the mass balance in particle accelerator experiments. This suggests that there must be a lower limit for the mass of black holes. Theoretically, this boundary is expected to lie around the Planck mass ($m_P = \sqrt{\hbar c/G} \approx 1.2 \times 10^{19}$ GeV/$c^2 \approx 2.2 \times 10^{-8}$ kg), where quantum effects are expected to make the theory of general relativity break down completely. This would put the creation of black holes firmly out of reach of any high energy process occurring on or near the Earth. Certain developments in quantum gravity however suggest that the Planck mass could be much lower: some braneworld scenarios

for example put it much lower, maybe even as low as 1 TeV/$c^2$ This would make it possible for micro black holes to be created in the high energy collisions occurring when cosmic rays hit the Earth's atmosphere, or possibly in the new Large Hadron Collider at CERN. These theories are however very speculative, and the creation of black holes in these processes is deemed unlikely by many specialists. Even if such micro black holes should be formed in these collisions, it is expected that they would evaporate in about $10^{-25}$ seconds, posing no threat to Earth

## Growth

Once a black hole has formed, it can continue to grow by absorbing additional matter. Any black hole will continually absorb gas and interstellar dust from its direct surroundings and omnipresent cosmic background radiation. This is the primary process through which supermassive black holes seem to have grown. A similar process has been suggested for the formation of intermediate-mass black holes in globular clusters.

Another possibility is for a black hole to merge with other objects such as stars or even other black holes. This is thought to have been important especially for the early development of supermassive black holes, which are thought to have formed from the coagulation of many smaller objects. The process has also been proposed as the origin of some intermediate-mass black holes.

## Evaporation

In 1974, Stephen Hawking showed that black holes are not entirely black but emit small amounts of thermal radiation. He got this result by applying quantum field theory in a static black hole background. The result of his calculations is that a black hole should emit particles in a perfect black body spectrum. This effect has become known as Hawking radiation. Since Hawking's result, many others have verified the effect through various methods. If his theory of black hole radiation is correct, then black holes are expected to emit a thermal spectrum of radiation, and thereby lose mass, because according to the theory of relativity mass is just highly condensed energy ($E = mc^2$). Black holes will shrink and evaporate over time. The temperature of this spectrum (Hawking temperature) is proportional to the surface gravity of the black hole, which for a Schwarzschild black hole is inversely proportional to the mass. Large black holes, therefore, emit less radiation than small black holes.

A stellar black hole of one solar mass has a Hawking temperature of about 100 nanokelvins. This is far less than the 2.7 K temperature of the cosmic microwave background. Stellar mass (and larger) black holes receive more mass from the cosmic microwave background than they emit through Hawking radiation and will thus grow instead of shrink. To have a Hawking temperature larger than 2.7 K (and be able to evaporate), a black hole needs to be lighter than the Moon (and therefore a diameter of less than a tenth of a millimeter).
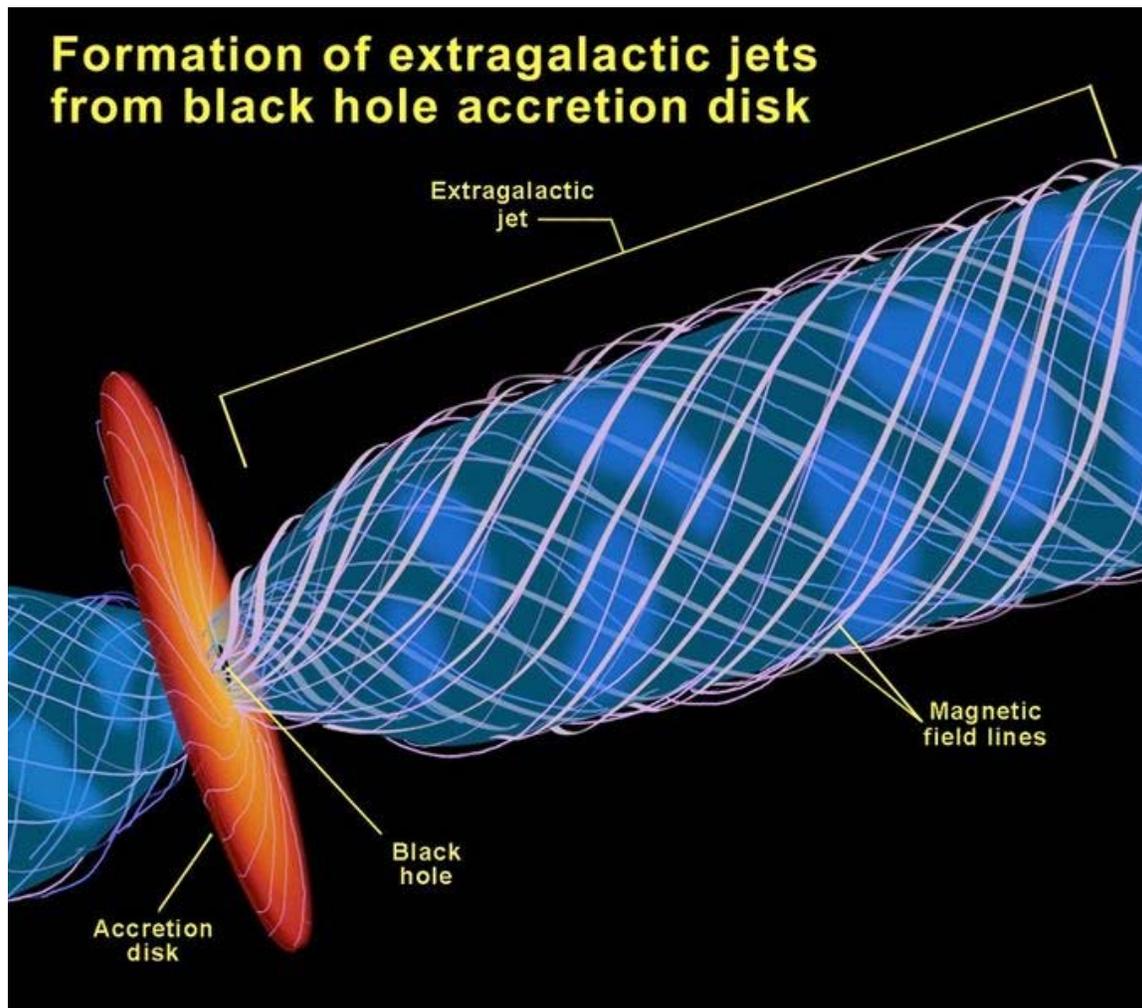
On the other hand, if a black hole is very small the radiation effects are expected to become very strong. Even a black hole that is heavy compared to a human would evaporate in an instant. A black hole the weight of a car ($\sim 10^{-24}$ m) would only take a nanosecond to evaporate, during which time it would briefly have a luminosity more than 200 times that of the sun. Lighter black holes are expected to evaporate even faster, for example a black hole of mass 1 TeV/$c^2$ would take less than $10^{-88}$ seconds to evaporate completely. Of course, for such a small black hole quantum gravitation effects are expected to play an important role and could even – although current developments in quantum gravity do not indicate so – hypothetically make such a small black hole stable.

## Observational evidence

By their very nature, black holes do not directly emit any signals other than the hypothetical Hawking radiation; since the Hawking radiation for an astrophysical black hole is predicted to be very weak, this makes it impossible to directly detect astrophysical black holes from the Earth. A possible exception to the Hawking radiation being weak is the last stage of the evaporation of light (primordial) black holes; searches for such flashes in the past has proven unsuccessful and provides stringent limits on the possibility of existence of light primordial black holes. NASA's Fermi Gamma-ray Space Telescope launched in 2008 will continue the search for these flashes.

Astrophysicists searching for black holes thus have to rely on indirect observations. A black hole's existence can sometimes be inferred by observing its gravitational interactions with its surroundings.
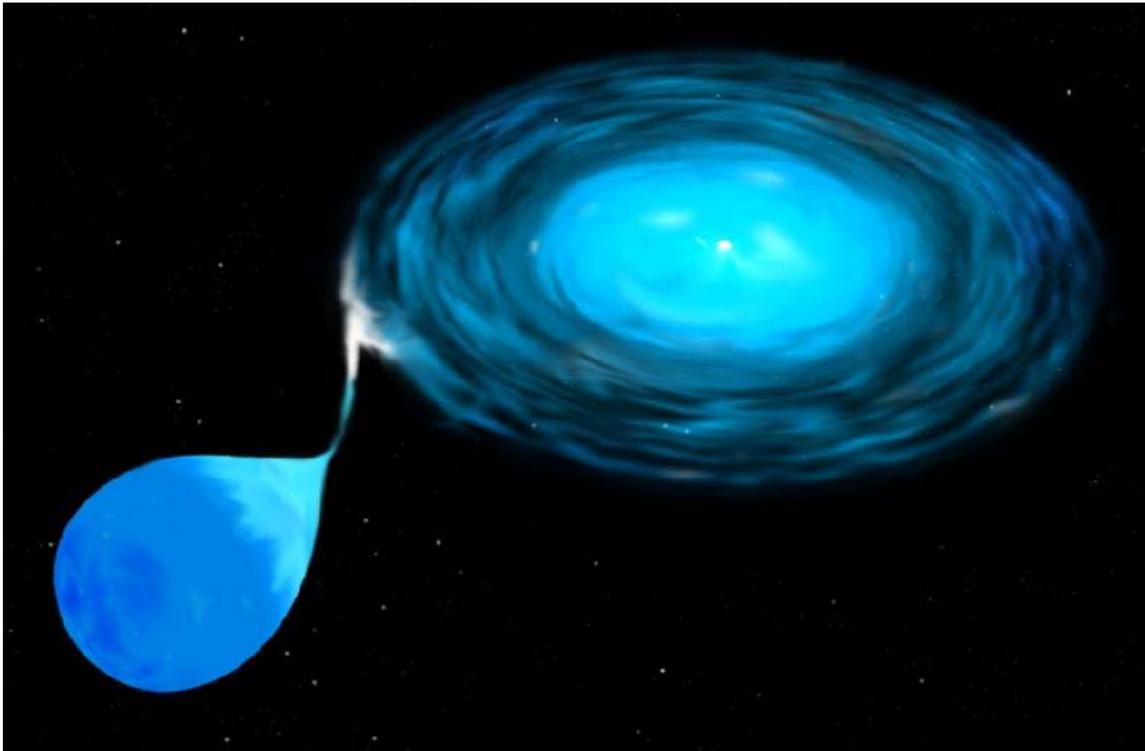
**Accretion of matter**



Formation of extragalactic jets from a black hole's accretion disk

Due to conservation of angular momentum, gas falling into the gravitational well created by a massive object will typically form a disc-like structure around the object. Friction within the disc causes angular momentum to be transported outward allowing matter to fall further inward releasing potential energy and increasing the temperature of the gas. In the case of compact objects such as white dwarfs, neutron stars, and black holes, the gas in the inner regions becomes so hot that it will emit vast amounts of radiation (mainly X-rays), which may be detected by telescopes. This process of accretion is one of the most efficient energy producing process known; up to 40% of the rest mass of the accreted material can be emitted in radiation. (In nuclear fusion only about 0.7% of the rest mass will be emitted as energy.) In many cases, accretion discs are accompanied by relativistic jets emitted along the poles, which carry away much of the energy. The mechanism for the creation of these jets is currently not well understood.

As such many of the universe's more energetic phenomena have been attributed to the accretion of matter on black holes. In particular, active galactic nuclei and quasars are thought to be the accretion discs of supermassive black holes. Similarly, X-ray binaries are thought to be binary star systems in which one of the two stars is a compact object accreting matter from its companion. It has also been suggested that some ultraluminous X-ray sources may be the accretion disks of intermediate-mass black holes.

## X-ray binaries

X-ray binaries are binary star systems that are luminous in the X-ray part of the spectrum. These X-ray emissions are generally thought to be caused by one of the component stars being a compact object accreting matter from the other (regular) star. The presence of an ordinary star in such a system provides a unique opportunity for studying the central object and determining if it might be a black hole.



Artist impression of a binary system with an accretion disk around a compact object being fed by material from the companion star.

If such a system emits signals that can be directly traced back to the compact object, it cannot be a black hole. The absence of such a signal does, however, not exclude the possibility that the compact object is a neutron star. By studying the companion star it is often possible to obtain the orbital parameters of the system and obtain an estimate for the mass of the compact object. If this is much larger than the Tolman–Oppenheimer–

Volkoff limit (that is, the maximum mass a neutron star can have before collapsing) then the object cannot be a neutron star and is generally expected to be a black hole.

The first strong candidate for a black hole, Cygnus X-1, was discovered in this way by Charles Thomas Bolton and Webster and Murdin in 1972. Some doubt, however, remained due to the uncertainties resultant from the companion star being much heavier than the candidate black hole. Currently, better candidates for black holes are found in a class of X-ray binaries called soft X-ray transients. In this class of system the companion star is relatively low mass allowing for more accurate estimates in the black hole mass. Moreover, these systems are only active in X-ray for several months once every 10–50 years. During the period of low X-ray emission (called quiescence), the accretion disc is extremely faint allowing for detailed observation of the companion star during this period. One of the best such candidates is V404 Cyg.

### Quiescence and advection-dominated accretion flow

The faintness of the accretion disc during quiescence is thought to be caused by the flow entering a mode called an advection-dominated accretion flow (ADAF). In this mode, almost all the energy generated by friction in the disc is swept along with the flow instead of radiated away. If this model is correct, then it forms strong qualitative evidence for the presence of an event horizon. Because, if the object at the center of the disc had a solid surface, it would emit large amounts of radiation as the highly energetic gas hits the surface, an effect that is observed for neutron stars in a similar state.
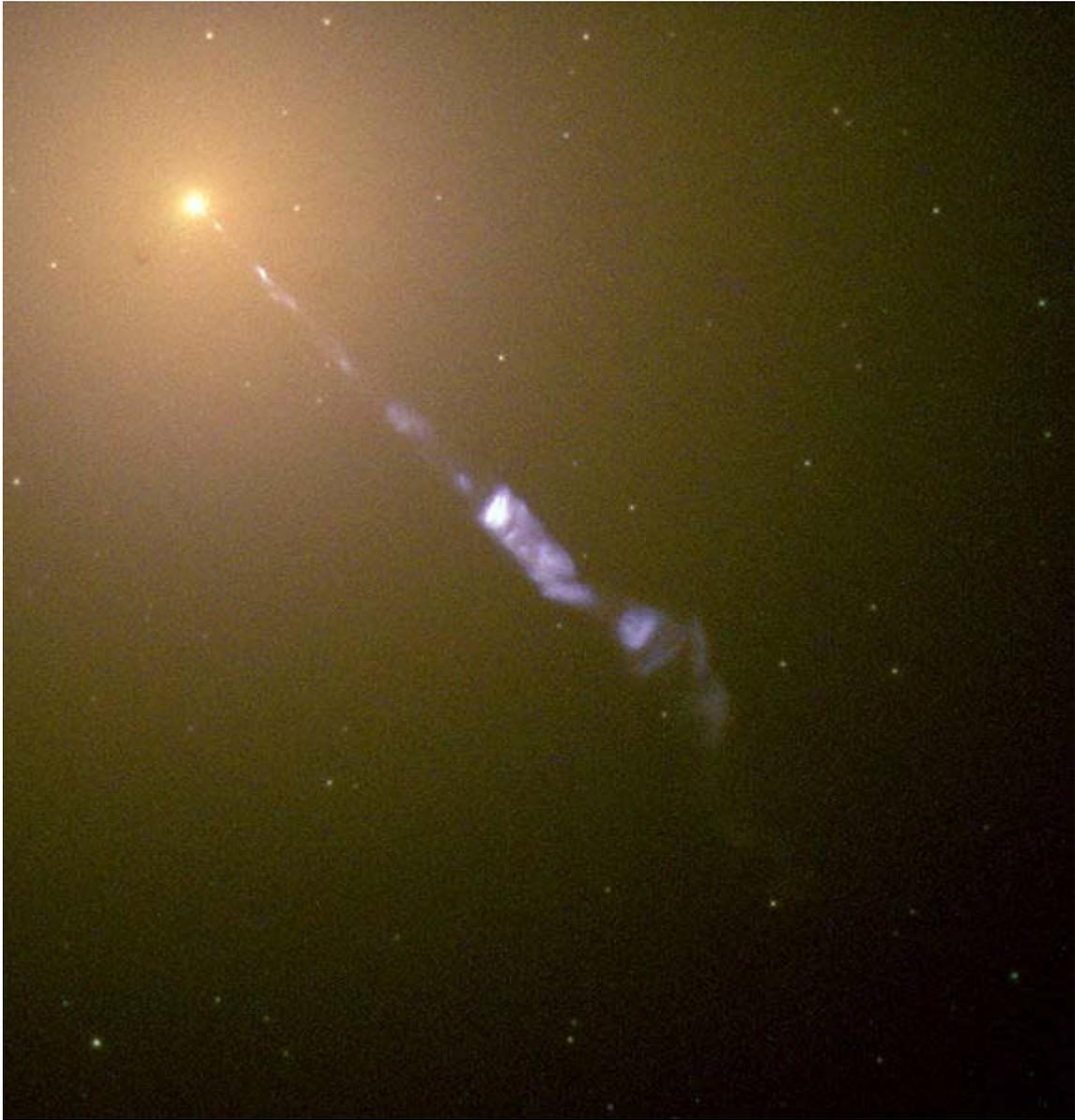
### Quasi-periodic oscillations

The X-ray emissions from accretion disks sometimes exhibit a flickering around certain frequencies. These signals are called quasi-periodic oscillations and are thought to be caused by material moving along the inner edge of the accretion disk (the innermost stable circular orbit). As such their frequency is linked to the mass of the compact object. They can thus be used as an alternative way to determine the mass of potential black holes.

## Gamma ray bursts

Intense but one-time gamma ray bursts (GRBs) may signal the birth of "new" black holes, because astrophysicists think that GRBs are caused either by the gravitational collapse of giant stars or by collisions between neutron stars, and both types of event involve sufficient mass and pressure to produce black holes. It appears that a collision between a neutron star and a black hole can also cause a GRB, so a GRB is not proof that a "new" black hole has been formed. All known GRBs come from outside our own galaxy, and most come from billions of light years away so the black holes associated with them are billions of years old.

**Galactic nuclei**



The jet originating from the center of M87 in this image comes from an active galactic nucleus that may contain a supermassive black hole. Credit: Hubble Space Telescope/NASA/ESA.

It is now widely accepted that the center of every or at least nearly every galaxy contains a supermassive black hole. The close observational correlation between the mass of this hole and the velocity dispersion of the host galaxy's bulge, known as the M-sigma relation, strongly suggests a connection between the formation of the black hole and the galaxy itself.

For decades, astronomers have used the term "active galaxy" to describe galaxies with unusual characteristics, such as unusual spectral line emission and very strong radio emission. However, theoretical and observational studies have shown that the active galactic nuclei (AGN) in these galaxies may contain supermassive black holes. The models of these AGN consist of a central black hole that may be millions or billions of times more massive than the Sun; a disk of gas and dust called an accretion disk; and two jets that are perpendicular to the accretion disk.

Although supermassive black holes are expected to be found in most AGN, only some galaxies' nuclei have been more carefully studied in attempts to both identify and measure the actual masses of the central supermassive black hole candidates. Some of the most notable galaxies with supermassive black hole candidates include the Andromeda Galaxy, M32, M87, NGC 3115, NGC 3377, NGC 4258, and the Sombrero Galaxy.

Currently, the best evidence for a supermassive black hole comes from studying the proper motion of stars near the center of our own Milky Way. Since 1995 astronomers have tracked the motion of 90 stars in a region called Sagittarius A*. By fitting their motion to Keplerian orbits they were able to infer in 1998 that 2.6 million solar masses must be contained in a volume with a radius of 0.02 lightyears. Since then one of the stars—called S2—has completed a full orbit. From the orbital data they were able to place better constraints on the mass and size of the object causing the orbital motion of stars in the Sagittarius A* region, finding that there is a spherical mass of 4.3 million solar masses contained within a radius of less than 0.002 lightyears. While this is more than 3000 times the Schwarzschild radius corresponding to that mass, it is at least consistent with the central object being a supermassive black hole, and no "realistic cluster [of stars] is physically tenable."

## Gravitational lensing

The deformation of spacetime around a massive object causes light rays to be deflected much like light passing through an optic lens. This phenomenon is known as gravitational lensing. Observations have been made of weak gravitational lensing, in which photons are deflected by only a few arcseconds. However, it has never been directly observed for a black hole. One possibility for observing gravitational lensing by a black hole would be to observe stars in orbit around the black hole. There are several candidates for such an observation in orbit around Sagittarius A*.

## Alternatives

The evidence for stellar black holes strongly relies on the existence of an upper limit for the mass of a neutron star. The size of this limit heavily depends on the assumptions made about the properties of dense matter. New exotic phases of matter could push up this bound. A phase of free quarks at high density might allow the existence of dense quark stars, and some supersymmetric models predict the existence of Q stars. Some extensions of the standard model posit the existence of preons as fundamental building blocks of quarks and leptons which could hypothetically form preon stars. These
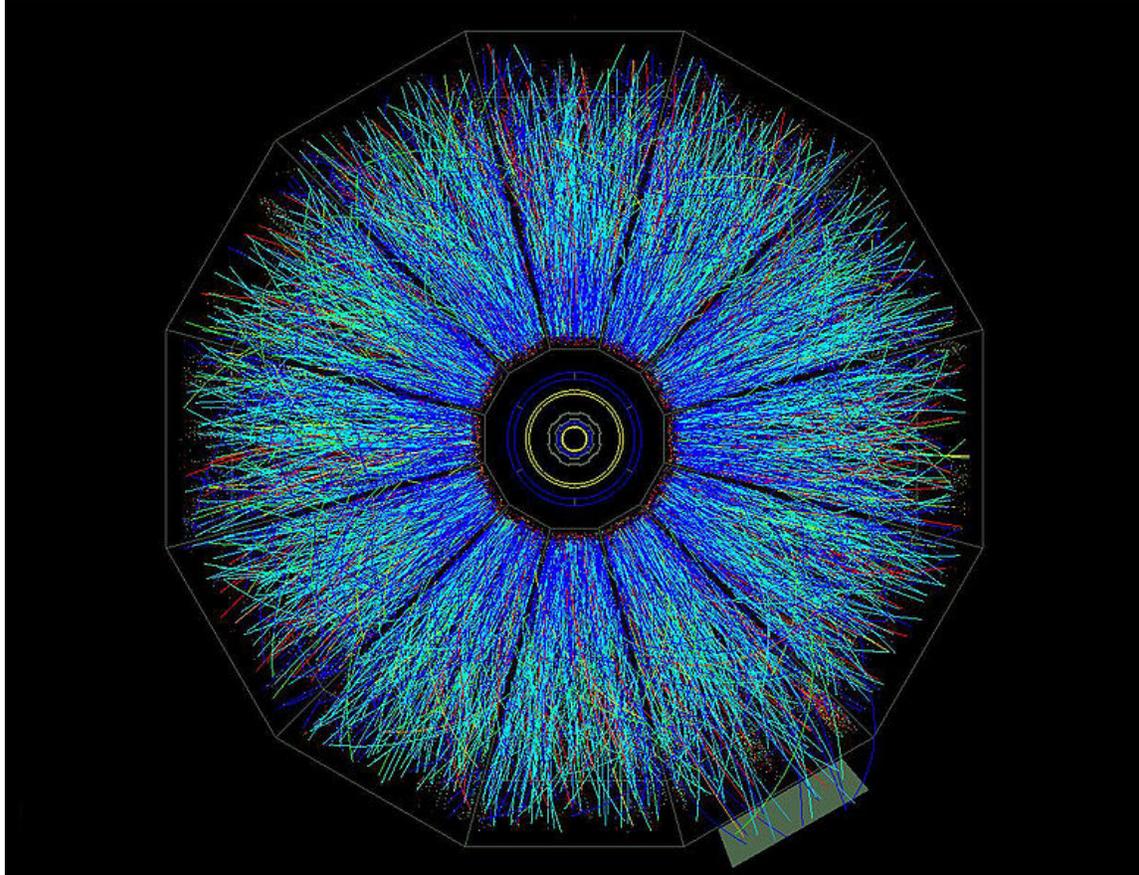
hypothetical models could potentially explain a number of observations of stellar black hole candidates. However, it can be shown from general arguments in general relativity that any such object will have a maximum mass.

Since the average density of a black hole inside its Schwarzschild radius is inversely proportional to the square of its mass, supermassive black holes are much less dense than stellar black holes (the average density of a large supermassive black hole is comparable to that of water). Consequently, the physics of matter forming a supermassive black hole is much better understood and the possible alternative explanations for supermassive black hole observations are much more mundane. For example, a supermassive black hole could be modelled by a large cluster of very dark objects. However, typically such alternatives are not stable enough to explain the supermassive black hole candidates.

The evidence for stellar and supermassive black holes implies that in order for black holes not to form, general relativity must fail as a theory of gravity, perhaps due to the onset of quantum mechanical corrections. A much anticipated feature of a theory of quantum gravity is that it will not feature singularities or event horizons (and thus no black holes). In recent years, much attention has been drawn by the fuzzball model in string theory. Based on calculations in specific situations in string theory, the proposal suggest that generically the individual states of a black hole solution do not have an event horizon or singularity (and can thus not really be considered to be a black hole), but that for a distant observer the statistical average of such states does appear just like an ordinary black hole in general relativity.

# Open questions

## Entropy and thermodynamics



If ultra-high-energy collisions of particles in a particle accelerator can create microscopic black holes, it is expected that all types of particles will be emitted by black hole evaporation, providing key evidence for any grand unified theory. Above are the high energy particles produced in a gold ion collision on the RHIC.

In 1971, Stephen Hawking showed under general conditions that the total area of the event horizons of any collection of classical black holes can never decrease, even if they collide and merge. This result, now known as the second law of black hole mechanics, is remarkably similar to the second law of thermodynamics, which states that the total entropy of a system can never decrease. As with classical objects at absolute zero temperature, it was assumed that black holes had zero entropy. If this were the case, the second law of thermodynamics would be violated by entropy-laden matter entering the black hole, resulting in a decrease of the total entropy of the universe. Therefore, Jacob Bekenstein proposed that a black hole should have an entropy, and that it should be proportional to its horizon area.

The link with the laws of thermodynamics was further strengthened by Hawking's discovery that quantum field theory predicts that a black hole radiates blackbody radiation at a constant temperature. This seemingly causes a violation of the second law of black hole mechanics, since the radiation will carry away energy from the black hole causing it to shrink. The radiation, however also carries away entropy, and it can be proven under general assumptions that the sum of the entropy of the matter surrounding the black hole and one quarter of the area of the horizon as measured in Planck units is in fact always increasing. This allows the formulation of the first law of black hole mechanics as an analogue of the first law of thermodynamics, with the mass acting as energy, the surface gravity as temperature and the area as entropy.

One puzzling feature is that the entropy of a black hole scales with its area rather than with its volume, since entropy is normally an extensive quantity that scales linearly with the volume of the system. This odd property led 't Hooft and Susskind to propose the holographic principle, which suggests that anything that happens in volume of spacetime can be described by data on the boundary of that volume.

Although general relativity can be used to perform a semi-classical calculation of black hole entropy, this situation is theoretically unsatisfying. In statistical mechanics, entropy is understood as counting the number of microscopic configurations of a system which have the same macroscopic qualities (such as mass, charge, pressure, etc.). Without a satisfactory theory of quantum gravity, one cannot perform such a computation for black holes. Some progress has been made in various approaches to quantum gravity. In 1995, Strominger and Vafa showed that counting the microstates of a specific supersymmetric black hole in string theory reproduced the Bekenstein–Hawking entropy. Since then, similar results have been reported for different black holes both in string theory and in other approaches to quantum gravity like loop quantum gravity.

## Black hole unitarity

An open question in fundamental physics is the so-called information loss paradox, or black hole unitarity paradox. Classically, the laws of physics are the same run forward or in reverse (T-symmetry). Liouville's theorem dictates conservation of phase space volume, which can be thought of as "conservation of information", so there is some problem even in classical physics. In quantum mechanics, this corresponds to a vital property called unitarity, which has to do with the conservation of probability (it can also be thought of as a conservation of quantum phase space volume as expressed by the density matrix).