# Theories of Consciousness
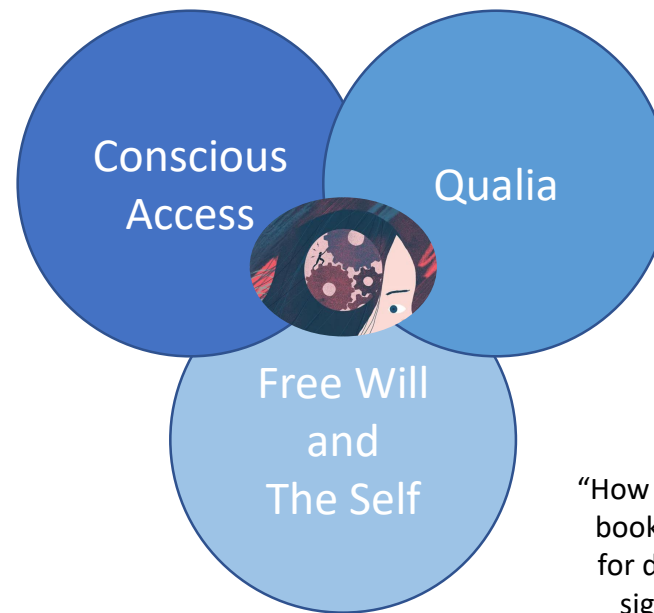
Speaker: Felix Haas

**Universität Zürich** UZH

# Definitions of Consciousness

- A theory of Consciousness is not going to give us a satisfactory definition of Consciousness.
- Rather, it starts with a definition of the phenomenon you want to describe.

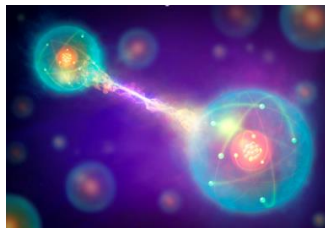"What counts as genuine consciousness,… is conscious access." Stanislas Dehaene in [Dehaene2014]

"The sensation of color cannot be accounted for by the physicist's objective picture of light-waves." Erwin Schrödinger in [Schroed1944]

Conscious Access

Qualia

Free Will and The Self

"How the Self Controls Its Brain," 1994 book by Sir John Eccles (Nobel Prize for discoveries on chemical basis of signal transmission at synapses)
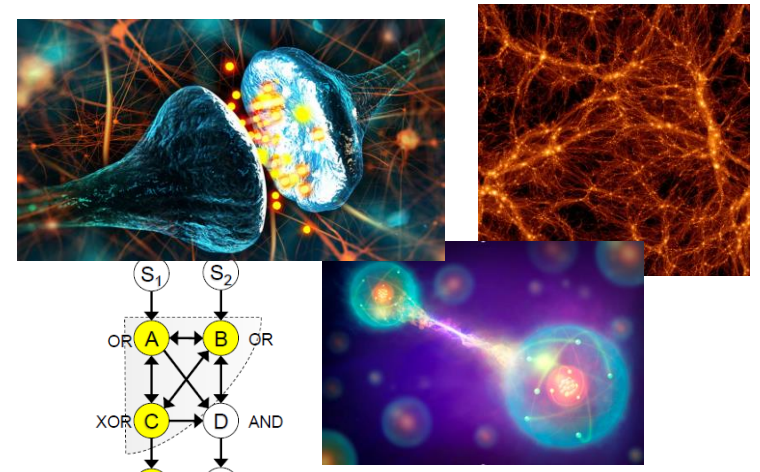
Universität Zürich UZH

# Levels of a future theory

Thermodynamics



Consciousness



Is this a fair comparison/ expectation?

What about the hard problem?



Universität Zürich UZH

# (Proto)theories of Consciousness for discussion

## Global Workspace



### Mathematical Formalism?

## Integrated Information



### Testability?

## Orchestrated Objective Reduction
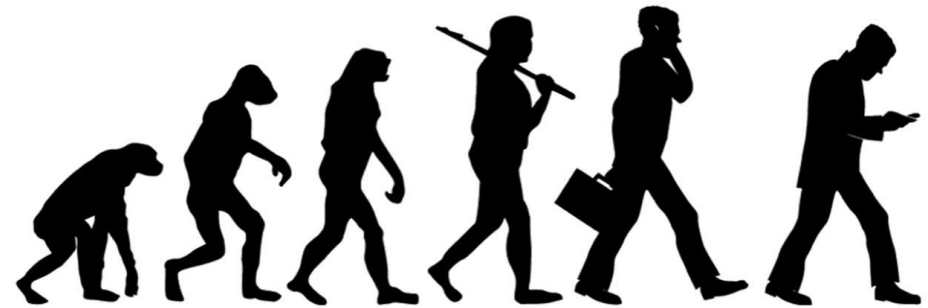


All treat Consciousness as emergent and (to a varying degree) as substrate independent.
What about Pan-psychism?

# Global (Neuronal) Workspace Theory (G(N)WT)

# Genesis of GWT – Psychological Insight

- Jerry Fodor's book, The Modularity of Mind [Fodor1983].

- Leda Cosmides and John Tooby lay modern foundations of **evolutionary psychology** in 1980s and 1990s: "…what emotions do… is to activate and coordinate the modular functions [e.g. jealousy, mating] that are, in Darwinian terms, appropriate for the moment." [Wright2017]

- "…that modules are triggered by feelings – sheds new light on the connection between two fundamental parts of **Buddhism**: the idea of non-attachments to feelings and the idea of not-self." [Wright2017]

[Fodor1983] Fodor, Jerry A., (1983), *The Modularity of Mind*, (MIT Press)
[Wright2017] Wright, Robert (2017), *Why Buddhism is True: The Science and Philosophy of Meditation and Enlightenment,* (Simon & Schuster)
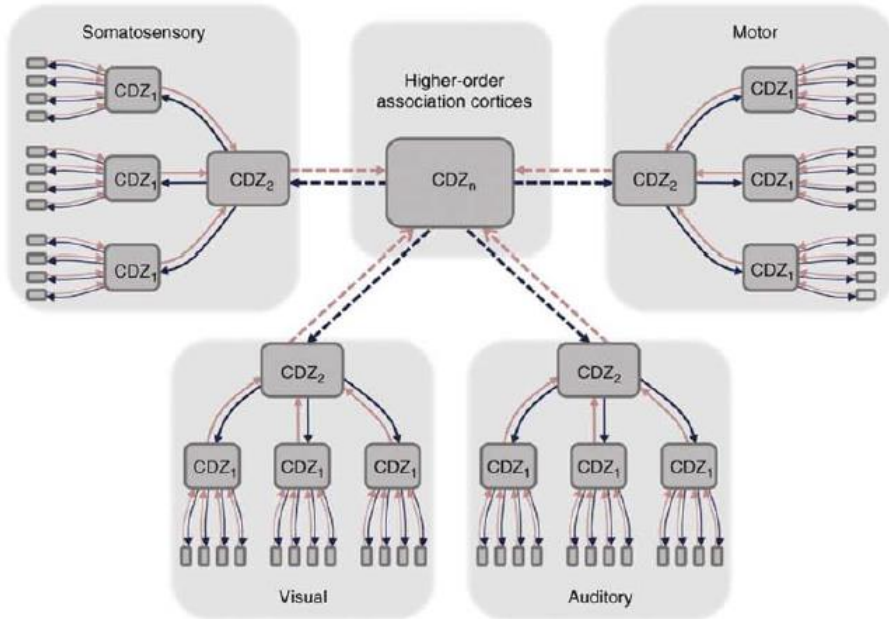
# Genesis of GWT – Neuroscience Insight



Image taken from [MKDD2013]

- Convergence-Divergence Zones (**CDZ**) [Dam1989], are **neural networks which function as specialized modules** with high connectivity to other parts of the brain.

- When a stimulus excites a specific set of the CDZ's neurons, this **strengthens the synapses connecting this set, forming a self-excitatory network**. A later excitation of this network is then capable of reproducing the original stimulus (**memory, imagination**).

- Cross modular divergence: Reading words of auditory or olfactory meaning (a purely visual input) activates specific networks in the auditory or olfactory cortical areas.

**Universität Zürich** UZH

# Global Workspace Theory (GWT)

- GWT was first put forward by Dutch psychologist Bernard Baars [Baars1988].

- Unconscious specialized modules compete for access to the Global Workspace, which integrates specialized information into a coherent interpretation of reality.

- Consciousness = Global broadcasting back to the specialized modules ("fame in the brain" [Dennett2005]).



**The Theater Metaphor**

Selective attention controls the spotlight that selects what will be in the bright spot on stage

Backstage is unconscious

The bright spot on stage has limited capacity.

The audience has vast capacity.

The audience is unconscious

From Carl Carpenter, *A New Model of Consciousnes*, Sci & Con Rev.2006.

- Spotlight = Consciousness
- Stage = Working memory
- Audience = Specialized unconscious processes (e.g. memory, language, automatisms)
- Backstage = intentions, expectations, self

[Baars1988] Baars, Bernard J. (1988), *A Cognitive Theory of Consciousness* (Cambridge University Press)
[Dennett2005] Dennett, Daniel, (2005), *Sweet Dreams: Philosophical Obstacles to a Science of Consciousness* (MIT Press)

**Universität Zürich** UZH

# Global Neuronal Workspace Theory (GNWT)

- Take Consciousness = Conscious Access.
- **Conscious states are encoded by the stable activation (i.e. for a few 100ms's) of a subset of active global workspace (GW) neurons**.

- The **GNW is a highly non-localized set of neurons** with long axions, connecting different unconscious modular neural networks.

- **Prefrontal cortex is "an important hub** of the global neuronal network … contributing to **non-linear ignition**." [MRCD2020]
- Entering the GW = becoming conscious = global "ignition" (of distributed brain regions) represents a **phase transition** in brain activity: Around **200 to 300 ms after stimulus onset**.
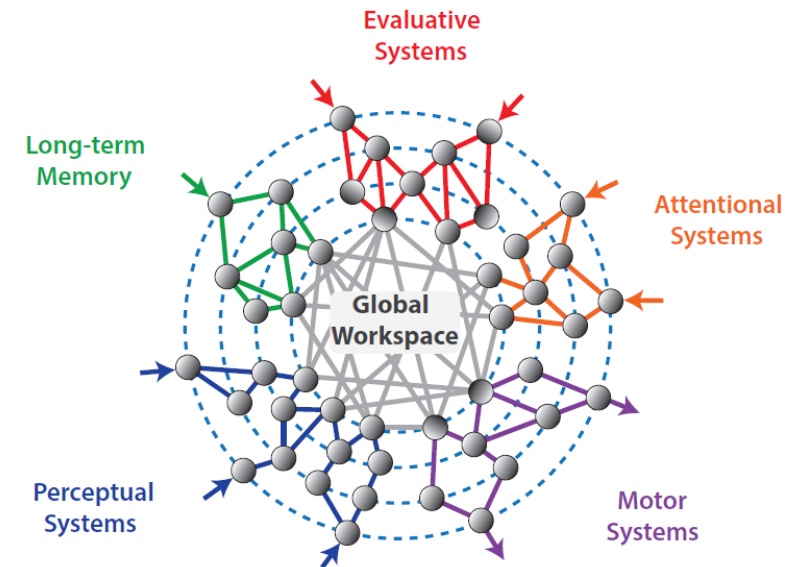


Bernard Baars          Stanislas Dehaene

**What is the GW anatomically?**

"Baars suggested the diffuse, extended **reticular-thalamic activating system** as the main brain structure forming the global workspace. However, Baars' instantiation of the hypothesis **does not distinguish between the level of conscious processing (under the control of the reticular formation) and the content**. By contrast, the GNW hypothesis, as initially proposed by Dehaene et al. (1998) and later simulated (Dehaene and Changeux, 2005; Dehaene et al., 2003), **proposes a defined brain network as the neural instantiation.**" [MRCD2020]
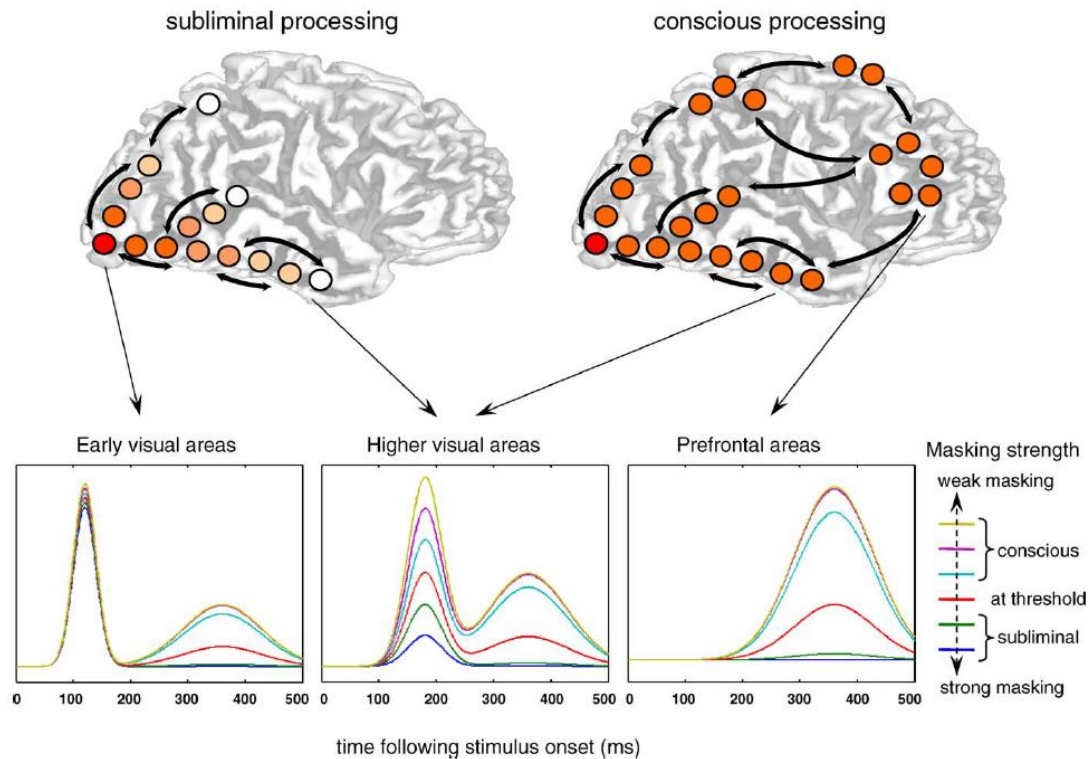
**Universität Zürich** UZH

# What GNWT helps us understand

- Why **Consciousness lags reality** by a quarter second.
- Why in **babies** up to 1 year old, **consciousness is even slower**: Distant cortical regions are already connected by long-distance fibers, but they are not yet covered in myelin (i.e. not yet well insulated).
- Why **masking** works: A picture remains subliminal ("below threshold") if flashed only for e.g. 33ms and if preceded and followed by masking images.
- **Schizophrenics**, who seem to struggle to integrate information into a coherent whole, show significant anomalies in their long-distance axons (particularly the corpus callosum), linking cortical regions.
- The **evolutionary purpose** of Consciousness is as a communication device between modular, local neuronal circuits.
- **Spontaneous neuronal activity** may help to push one thought/signal into Consciousness even though it was just subliminal.

# GNWT: Brain activation from masked stimuli

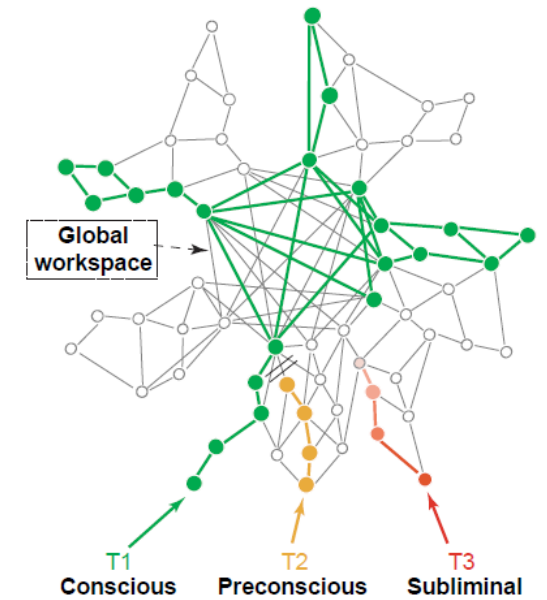GNWT schematic predition

Signatures of Consciousness



▶ Intense ignition of distributed brain regions, incl. bilateral prefrontal cortex.

▶ P300 (P3) wave

[DCBD2007] Del Cul A, Baillet S, Dehaene S (2007) *Brain dynamics underlying the nonlinear threshold for access to consciousness.* PLoS Biol 5(10): e260. doi:10.1371/journal.pbio.0050260
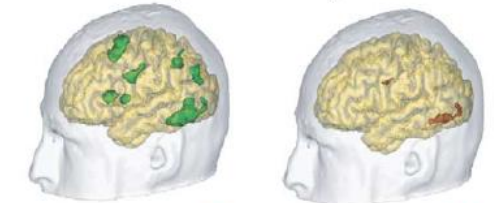
# GNWT – The Unconscious



According to GNWT, most brain activity is unconscious. GNWT helps us understand different types of unconsciousness:

- **Subliminal**: (e.g. incoming) sensory signal is too weak to ever become conscious.
- **Preconscious**: The signal is strong enough to become conscious but is not attended and so not further amplified to make it into the workspace.
- **Disconnected patterns**: Processors that are not connected to the workspace. E.g. respiration data forever remains in brain stem neurons.
- **Diluted information**: An individual signal pattern is diluted. E.g. differently oriented black-white gratings flicker so fast (>50 Hz), that you can only see gray, even though distinct visual neurons for the different gratings still fire.

Universität Zürich UZH

# GNWT – Road ahead
# (a selection)

- Further, more **detailed study of ignition** process and the brain regions' detailed roles in it.

- Further develop **computer simulations** that mirror the GNW and are to achieve relevant signatures of Consciousness.

- Deploy **GNWT insights to treatment**. E.g. jump-starting the GNW to help vegetative state patients to regain Consciousness.

- **Mathematical Formulation** of GNWT?

- **Qualia**?

Universität
Zürich UZH

For a bit more on GNWT, see:
Mashour, G., Roelfsema, P., Changeux, J.-P., Dehaene, S. (2020) Conscious Processing and the Global Neuronal Workspace Hypothesis. Neuron, Volume 105, Issue 5, 4 March 2020, Pages 776-798

# For (a lot) more on GNWT, see:

Also, Baars' original book is still a good read:
Baars, Bernard J. (1988), *A Cognitive Theory of Consciousness* (Cambridge University Press)

"A revealing and definitely not dumbed-down overview of what we know about consciousness." —*Kirkus Reviews*

# CONSCIOUSNESS AND THE BRAIN

DECIPHERING HOW THE BRAIN CODES OUR THOUGHTS

## STANISLAS DEHAENE

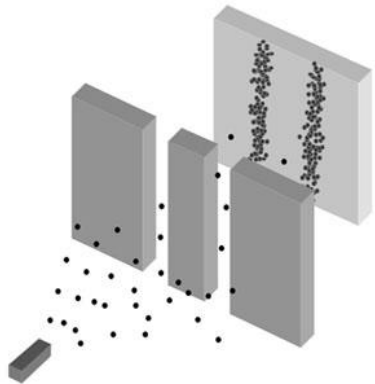Author of *The Number Sense* and *Reading in the Brain*

# Quantum Theory and Consciousness
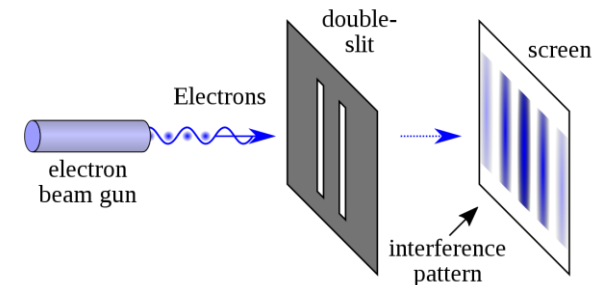
# What is quantum physics again?

**Classical Physics**
- Observables = functions of states (x,p)
- States and observables are perfectly sharp (cats are either 100% dead or 100% alive).
- Measurement is «outside of the theory» and in no way interferes with the state.

**Quantum Physics**
- Divides world into (quantum) "system" and "the rest". The system's states $| \Psi >$ are not directly observable and probabilistic.
- Possible superposition of states.
- Observables (things we want to measure, e.g. energy, position) are Operators on States.
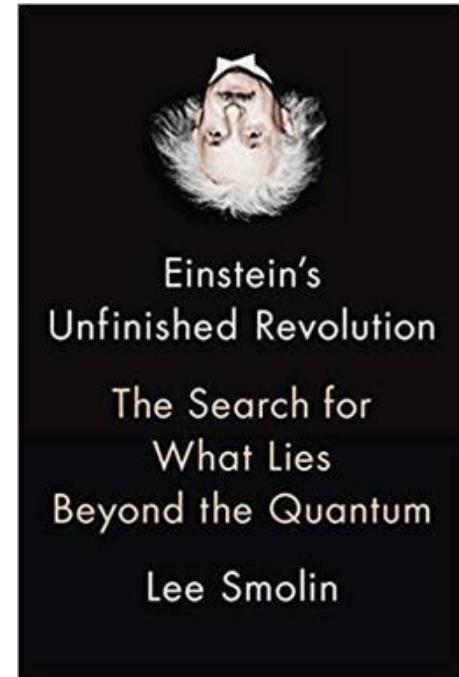- Measurement of an Observable will change the state.



- Discreteness of some quantities for some systems
- Heisenberg uncertainty
- Wave–particle duality
- Entanglement ("spooky action at a distance")

**Universität Zürich**UZH

16

# Quantum theory = unfinished business?

Issues with Quantum (Field) Theory (Q(F)T):

- **QT's axioms are intrinsically incompatible (Measurement Problem)*.**
- QT's axioms divide the universe into system and not-system.
- There are abundant mathematical inconsistencies in QFT.
- QFT cannot (to date) be solved analytically.

- QT is not a realist theory. Observables are only well defined when measured.
- Incompatible with General Relativity.

Einstein's
Unfinished Revolution

The Search for
What Lies
Beyond the Quantum

Lee Smolin

*In QM, a particle's evolution in time is governed by the (linear) Schrödinger equation (one axiom of QM), except for when it is measured, then the particle collapses (non-linearly) to the measurement outcome.

**Universität Zürich** UZH

# No quantum psychokinesis



- **The quantum measurement process altering the state of the system observed is sometimes taken to suggest that our consciousness collapses the wave function.**

- Even though it talks of "observables" and "observations", **quantum mechanics does not require a mind or consciousness**. In the double-slit experiment, the 2nd screen shows the pattern, no matter if a consciousness is watching.

- Yet still, you find claims like: "Probably the crux of quantum science is the relationship between consciousness and reality. ...Ultimately, panpsychism is grounded in, or is supported by, quantum entanglement." [M2018]

**Universität Zürich** UZH

# Why some might look to quantum physics to explain consciousness

- **Quantum indeterminism** in an otherwise fully deterministic universe, might be **viewed to open a door to free will**.
- **Quantum non-locality** somehow allows "microtubules in our brains [to act] like **antennae for consciousness**."
- Perhaps a longing for the unification of spirituality perceived to be inherent to quantum physics with our own.



**Many phycists' reactions** to being asked about quantum physics and Consciousness is an eye-roll, usually followed by one or both of the following arguments:

- **Copenhagen interpretation does not include or necessitate the presence of or interaction with a Consciousness.**
- **The "warm, wet and noisy" argument**.
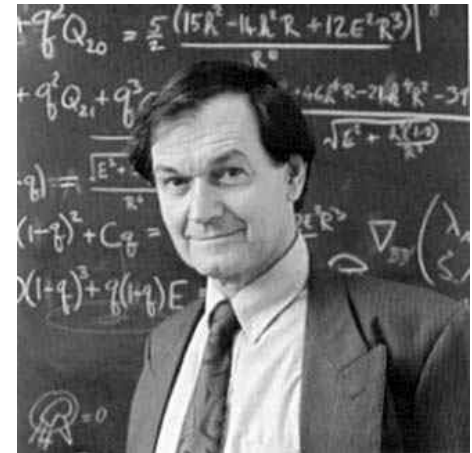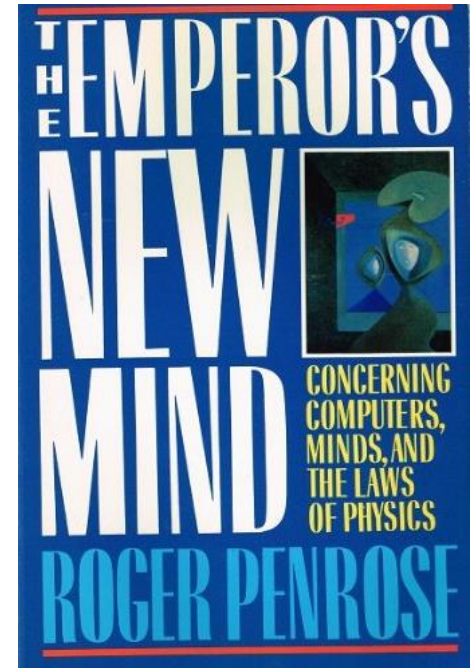  (See e.g. Lawrence Krauss: https://www.nbcnews.com/science/how-spot-quantum-quackery-6c10403763)

# Orchestrated Objective Reduction (Orch OR)

# Roger Penrose &
# The emperor's new mind



Penrose received the 2020 Nobel Prize in Physics "for the discovery that black hole formation is a robust prediction of the general theory of relativity."*

Penrose published "**The Emperor's New Mind**" in 1989, a project which he had started after hearing "extreme AI opinions" expressed by "proponents of strong AI." [Pen1989]

He deploys the "Penrose-Lucas** argument" to conclude that the **human mind could never be fully replicated or surpassed by a machine**.
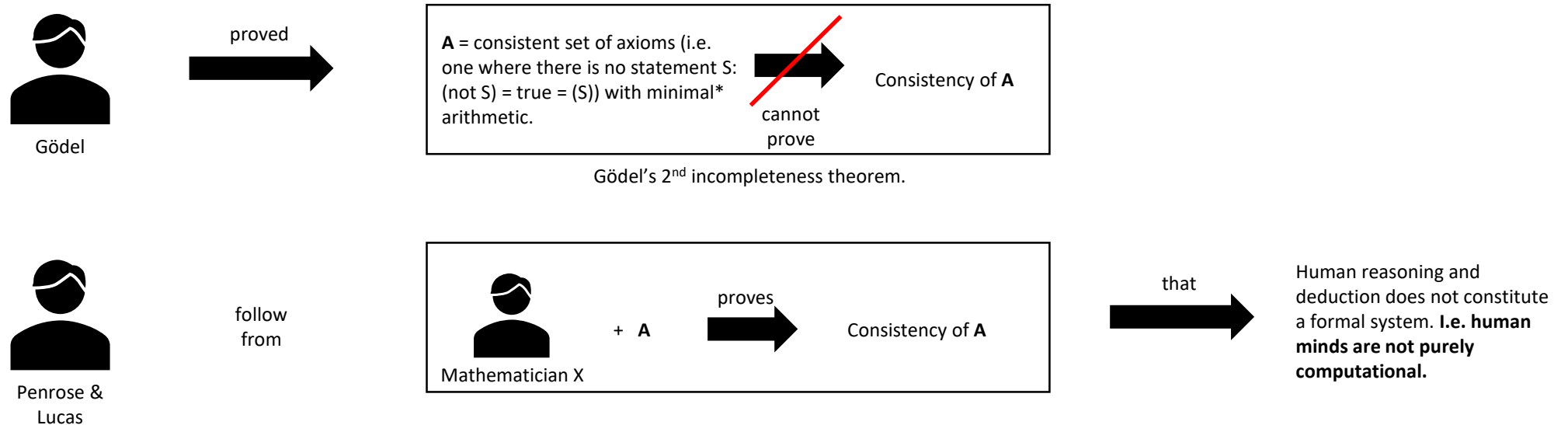


*Penrose–Hawking singularity theorems (1965 Penrose publication and 1970 publication of Penrose and Hawking). Outside of GR, he also invented Twistor theory and made significant contributions to the study of tessellations.
** Philosopher J.R. Lucas had made the same argument in 1961.
[Pen1989] Penrose, R. (1989), The Emperor's New Mind (Oxford University Press)

# Penrose-Lucas argument

**Gödel**

proved →

A = consistent set of axioms (i.e. one where there is no statement S: (not S) = true = (S)) with minimal* arithmetic.

→ Consistency of **A**

cannot prove

Gödel's 2nd incompleteness theorem.

**Penrose & Lucas**

follow from

Mathematician X + **A** proves → Consistency of **A**

that →

Human reasoning and deduction does not constitute a formal system. **I.e. human minds are not purely computational.**

There is a **wide number of criticisms and rejections** from computer scientists, mathematicians, philosophers. One argument made is that humans are not (or cannot be proven to be) consistent formal systems (see e.g. computer scientist M.L. Minsky). For a much more in-depth and technical criticism see e.g.: [Franzen2005], or [LHK1998].

*minimal arithmetic = (+,·, the symbols ∀,∃, and the usual rules for handling them)
[Franzen2005] Franzén, Torkel (2005). Gödel's Theorem: An Incomplete Guide to its Use and Abuse. (Wellesley, Massachusetts: A K Peters, Ltd.)
[LHK1998] LaForte, G., Hayes, P.J., Ford, K.M., (1998), Why Gödel's Theorem Cannot Refute Computationalism. Artificial Intelligence, 104:265–286.
Also: https://chronos-tachyon.net/essays/penrose-objections/

**Universität Zürich**UZH

# Penrose's inferences from the supposed non-computable nature of the human mind
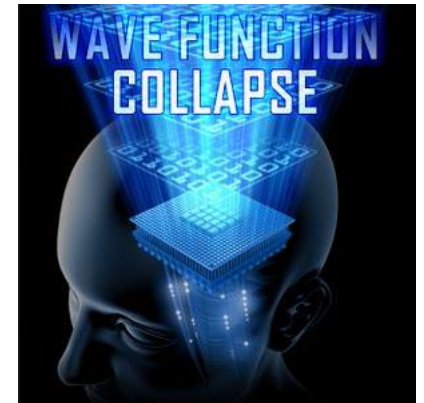
«Instead of Consciousness causing collapse, collapse causes [Proto-]Consciousness.» [Hameroff2019]

If the **quantum mechanical wave function collapse is** truly random, then this is a prime example of a **non-computable process**.

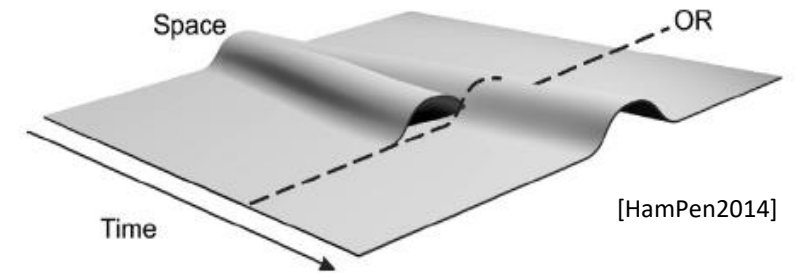- **Hypothesis 1 = Wave function collapse gives rise to a proto-consciousness**.

In the mainstream (Copenhagen) interpretation of quantum theory, wave function collapse only happens at measurement.

- **Hypothesis 2 = Objective Reduction** = The wave function collapse is a real physical process, not exclusive to measurements. Objective Reduction (or collapse) of the wave function, when a space-time curvature threshold is surpassed.

**Universität Zürich**UZH

# What is Objective Reduction (OR)?

- Penrose argues [Penrose1996] that:
  - Due to General Relativity, these masses or energies give rise to **2 distinct (superposed) Space-Time geometries**.
  - Following the Heisenberg Uncertainty Principle, this **indeterminacy becomes "macroscopically relevant,"** at $T \simeq \hbar/\Delta E$, and the superposed geometries collapse to just one.

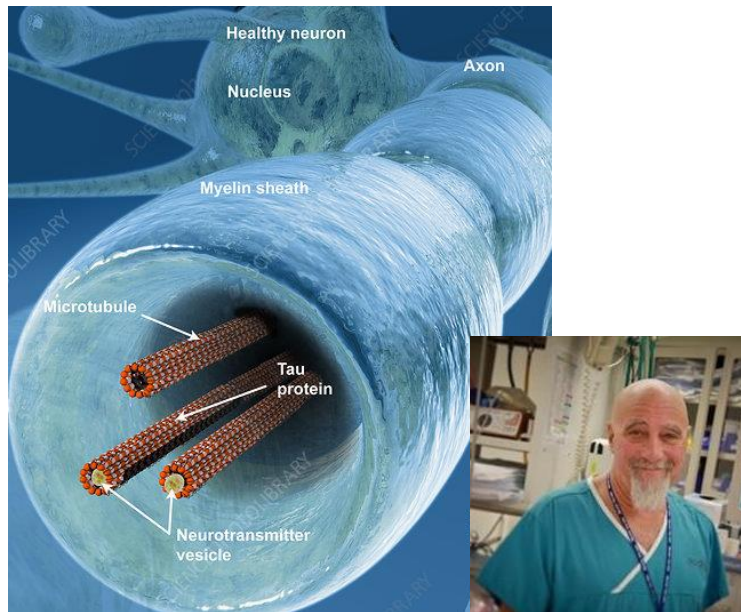**OR is not a generally accepted extension of quantum theory**.
- We have no theory of quantum gravity. Penrose effectively uses Schrödinger–Newton equation.
- OR violates unitarity and thus Energy conservation.



[HamPen2014]

$\Delta E$ is the self-energy of the difference between the two mass distributions.
[Penrose1996] R. Penrose, "On Gravity's Role in Quantum State Reduction," General Relativity and Gravitation 28: 581-600
[HamPen2014] S. Hameroff, R. Penrose, "Consciousness in the universe: A review of the 'Orch OR' theory," Physics of Life Reviews 11 (2014) 39-78:
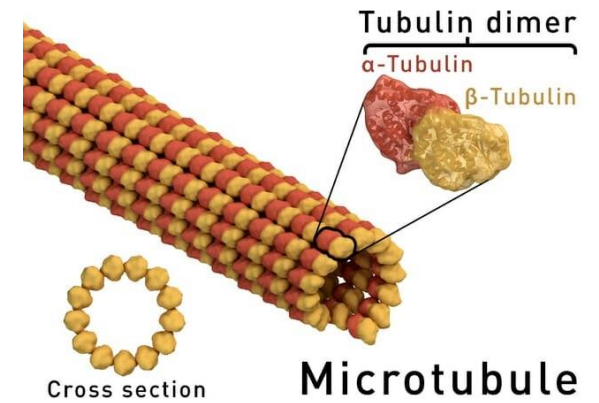https://www.sciencedirect.com/science/article/pii/S1571064513001188

# Orchestration – How does this pan-proto-psychism form our Consciousness?



- **Penrose needed a quantum computer** in the brain that **Orchestrates** these moments of Proto-Consciousness coming from Objective (collapse or) Reduction into the Consciousness we experience. **Hameroff suggested microtubules**.

- His idea was that **individual tubulin molecules** in the microtubules of a neuron **acts as** a quantum-bit, or "**qubit**," which are the elementary building blocks of quantum computing. Requiring the tubulin to be able to switch between alternative states in a coherent manner.
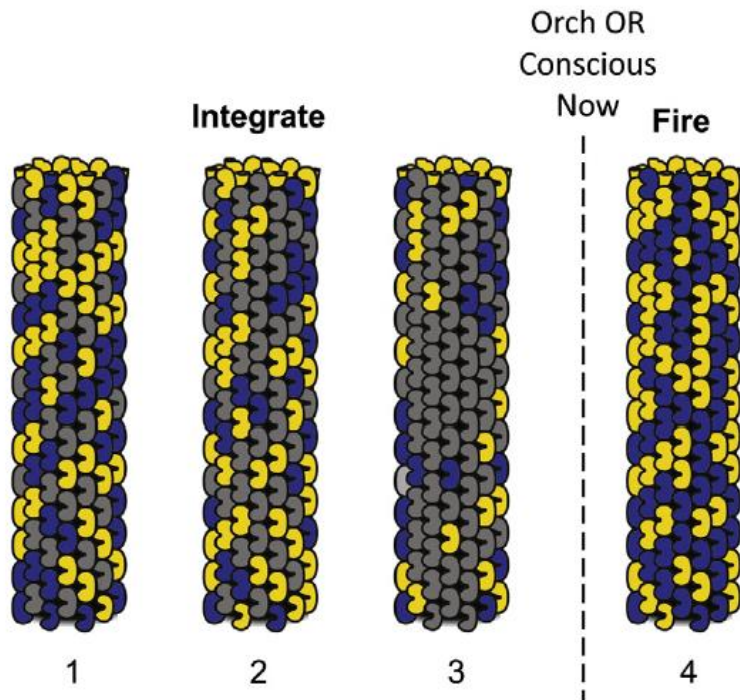
# Microtubules

- **Microtubules are found in all eukaryotic cells** — from humans to worms to sponges.
- They are built up from α and β tubulin proteins, each made from ca. 450 amino acids.

- Their known functions include:
  - **Integral part of the cytoskeleton** (the scaffolding of the cell) giving it shape and flexibility.
  - Important **role in cell division**, forming the "mitotic spindle" which organizes the chromosomes, pulls them apart, and steers them into the daughter cells' nucleus.
  - As a **highway for the transport of motor protein. E.g. axonal transport on microtubules in neurons:** motor proteins kinesin and dynein transport cargoes including mitochondria, cytoskeletal polymers, and synaptic vesicles containing neurotransmitters away from or towards the cell body.



Kinesin protein takes a walk on a microtubule:
https://www.youtube.com/watch?v=xlPDEpimzB8
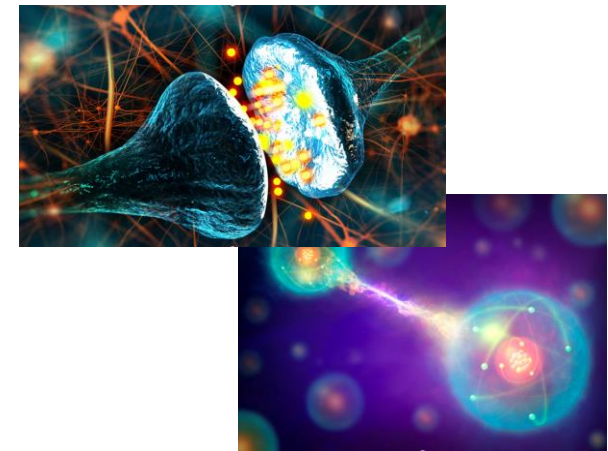
Universität Zürich UZH

# Orch OR – Conscious Now

S. Hameroff, R. Penrose / Physics of Life Reviews 11 (2014) 39–78



- "Tubulins are in classical **dipol states (yellow and blue),** or quantum **superposition of both states (grey)**. Quantum superposition/computation evolves during integration phases (1-3)… increasing quantum superposition… until threshold is met at time $T = \hbar/\Delta E$, at which time a conscious moment occurs…" [HamPen2014]

- "Tubulin dipoles in Orch OR were originally described in terms of **London force electric dipoles**, involving charge separation. However we now suggest, as an alternative, **magnetic dipoles**, which could be **related to electron spin** – and possibly **related also to nuclear spins**." [HamPen2014]

[HamPen2014] S. Hameroff, R. Penrose, "Consciousness in the universe: A review of the 'Orch OR' theory," Physics of Life Reviews 11 (2014) 39-78: https://www.sciencedirect.com/science/article/pii/S1571064513001188

Universität Zürich UZH

# Orch OR – Criticism

In addition to criticism of the Penrose-Lucas argument and
Objective Reduction (OR) already mentioned.



High-level criticism:
- No evidence that microtubules play a functional role in the signaling behavior of neurons.
- Significant evidence that neurochemical and electrical signaling play a central role in all aspects of brain function.
- If microtubules can function as quantum computers, why has evolution only used this power in brains and other neuronal structures?

Lower-level criticism:
- The "**warm, wet and noisy**" argument: Quantum computation requires isolation to prevent thermal interactions and decoherence, destroying quantum states [Tegmark2000]. Warm quantum coherence confirmed in plant photosynthesis or bird brain navigation does not imply that Orch OR is true.
- "Assuming that microtubule quantum states occur in a specific brain neuron, **how could it involve microtubules in other neurons throughout the brain?** Orch OR proposes that quantum states can extend by **entanglement between adjacent neurons through gap junctions**, primitive electrical connections between adjacent cells." [HamPen2014]
  A mechanism that has been deemed unlikely.
- How, in a controlled manner, does **classical information get codified into/later re-codified from** (entangled) dipol-qbits? How does the **entanglement preparation** happen in a controlled matter? What are the **quantum (logical) gates**? Etc.

[Tegmark2000] Tegmark, Max (2000). "Importance of quantum decoherence in brain processes". Physical Review E. 61 (4): 4194–4206. arXiv:quant-ph/9907009
[MRMH2009] McKemmish LK, Reimers JR, McKenzie RH, Mark AE, Hush NS. "Penrose-Hameroff orchestrated objective reduction proposal for human consciousness is not biologically feasible." Phys Rev E 2009; 80 (2Pt1): 021912
[HamPen2014] S. Hameroff, R. Penrose, "Consciousness in the universe: A review of the 'Orch OR' theory," Physics of Life Reviews 11 (2014) 39-78: https://www.sciencedirect.com/science/article/pii/S1571064513001188

Universität
Zürich UZH

# Orchestrated Objective Reduction (Orch OR) – Further Reading

- https://en.wikipedia.org/wiki/Orchestrated_objective_reduction
- www.quantumconsciousness.org
- S. Hameroff, R. Penrose, "Consciousness in the universe: A review of the 'Orch OR' theory," Physics of Life Reviews 11 (2014) 39-78: https://www.sciencedirect.com/science/article/pii/S1571064513001188
- [Hameroff2019] The Science of Consciousness: Stuart Hameroff: https://www.youtube.com/watch?v=JHg-mr4aqWk

- For an alternative angle on quantum effects in cognition, see the appendix, or directly: https://www.quantamagazine.org/a-new-spin-on-the-quantum-brain-20161102/

**Universität Zürich** UZH

# Appendix

Universität
Zürich UZH

# A Selection of (Proto)Theories of Consciousness

| | Quantum (Orch OR) | Global (Neuronal) Workspace | Integrated Information |
|---|---|---|---|
| **Consciousness is Emergent?** | Yes | Yes | Yes |
| **Consc. is substrate Independent?** | Yes, but need Orchestrator. | Yes in principle, but theory very much linked to brain. | Yes |
| **Pan-Psychism?** | Pan-Proto-Psychist | Not explicitly (possibly "compatible" with a version of Pan-Psychism). | Not in the classical sense, but even small logical circuits would have Consc. |
| **Mathematical Formalism?** | No | No | Yes |
| **Reductionism?** | Yes | Down to some brain regions, but not fully. | Yes |
| **Experimentally Falsifiable in current version?** | Elements of it should be, others probably not. | Most developed in terms of linking to observable NCC's and subjective reports. | Not at this stage. Discussions on measurable correlate metrics ongoing. |

| Meditation / Introspection |
|---|
| [Fill out] |

Universität
Zürich UZH

# Reticular Activating System (RAS)



- "The reticular activating system spans an extensive portion of the brainstem. [Its] fundamental role is regulating arousal and sleep–wake transitions. The ascending reticular activating system projects to the intralaminar nuclei of the thalami, which projects diffusely to the cerebral cortex. The ascending projections of the reticular activating system enhance the attentive state of the cortex and facilitate conscious perception of sensory stimuli." [WS2014]

- The RAS helps regulate what enters the global workspace (i.e. what becomes conscious), but itself is not the GW.

[WS2014] B.L. Walter, A.G. Shaikh, in Encyclopedia of the Neurological Sciences (Second Edition), 2014

# Signatures of Consciousness

Science of Consciousness in the late 80s and 90s:

1. Clearer (more limited) definition of "Consciousness," i.e. conscious access vs. qualia,
2. Consciousness manipulated (from masking and subliminal priming to deep brain stimulations),
3. Taking subjective reports seriously as data,
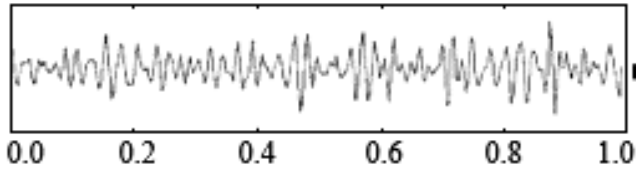4. Invention of functional MRI (fMRI) in 1990.

Signatures or Markers of Consciousness:

1. Intense ignition of distributed brain regions, incl. bilateral prefrontal cortex.
2. P300 (P3) wave (event-related potential (ERP) measured through EEG),
3. Late amplification (not just its mere presence) of gamma-band (= 30-100+ Hz) EEG-activity.
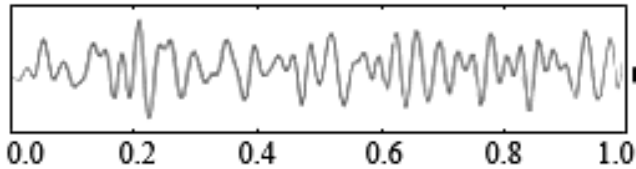4. Synchronization of electromagnetic signals across the cortex.

# EEG – Brain Waves
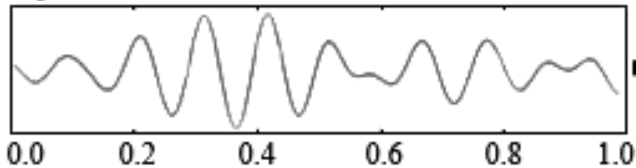


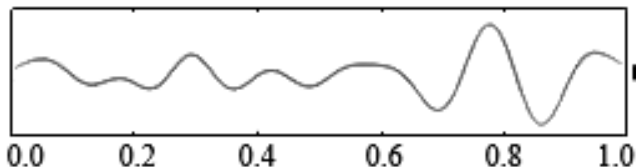Gamma Waves (30Hz-100Hz) → Motor Functions, higher mental

Beta Waves (12Hz-30Hz) → Normal waking state, concentration, focus, five physical sense, integrated
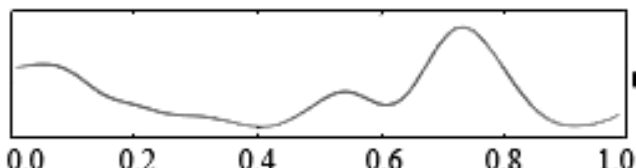
Alpha Waves (7.5Hz-12Hz) → Relaxed, light meditation, creative, super learning, conscious
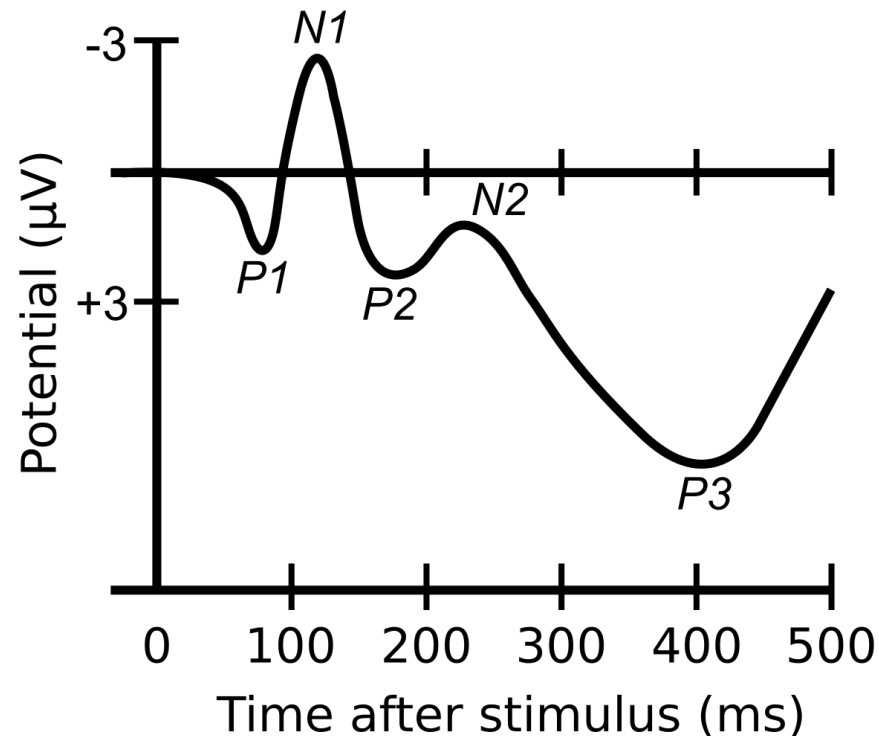
Theta Waves (4Hz-7.5Hz) → Light sleep, deep meditation, creative, recall, fantasy,

Delta Waves (0.1-4Hz) → Deep, dreamless sleep, non-REM sleep, unconscious

Description

**Universität Zürich** UZH

# Event-Related Potential (ERP)



- An ERP is the brain's measured response to a defined sensory, motor, or cognitive event.

- ERPs are measured by means of electroencephalography (EEG). The magnetoencephalography (MEG) equivalent of ERP is the ERF, or event-related field.

- Due to a multitude of simultaneous brain activity, you will not be able to see a clear ERP in single trials. Only when averaging over many results, will random brain activity be averaged out, with the remaining wave being referred to as ERP.

**Universität Zürich**<sup>UZH</sup>

# Kurt Gödel's 1931 two incompleteness theorems

Definitions

- **A consistent formal system** is any set of axioms which does not include any statement such that both the statement and its negation are provable from the axioms.

- A **set of axioms is complete** if, for any statement in the axioms' language, that statement or its negation is provable from the axioms.

- A **Gödel sentence G(F)** is a statement within F that is true iff G(F) cannot be proved in F – that is, it can be rendered in English as "this sentence is not provable in F".

Theorems

- **First Theorem**: In any consistent formal system F within which a certain amount of arithmetic (+,·, the symbols ∀,∃, and the usual rules for handling them) can be carried out, there are statements of the language of F which can neither be proved nor disproved in F.
  I.e. **All consistent formal systems with minimal arithmetic are incomplete**.

- **Second Theorem**: Any such formal system F cannot prove that the system itself is consistent (assuming it is indeed consistent).
  I.e. **A consistent formal system with minimal arithmetic cannot prove its consistency**.

An Example

- **An example of an undecidable problem is the continuum hypothesis** (= first of Hilbert's problems), advanced by Georg Cantor in 1878: Any infinite subset of the real numbers bijects either to the integers or the real numbers.
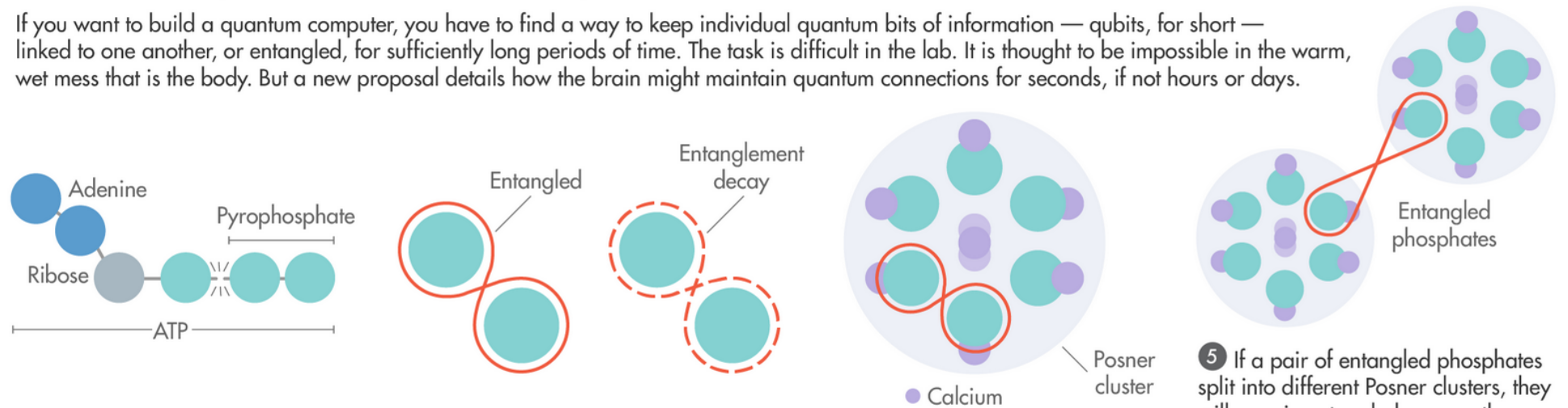
**Universität Zürich**UZH

# An example of an idea for "quantum cognition" with a smaller scope than Orch OR.

Matthew Fisher

"**Quantum Cognition: The possibility of processing with nuclear spins in the brains**," M.P.A. Fisher, Annals of Physics 362, p. 593-602 (2015):
https://www.kitp.ucsb.edu/sites/default/files/users/mpaf/174.pdf

If you want to build a quantum computer, you have to find a way to keep individual quantum bits of information — qubits, for short — linked to one another, or entangled, for sufficiently long periods of time. The task is difficult in the lab. It is thought to be impossible in the warm, wet mess that is the body. But a new proposal details how the brain might maintain quantum connections for seconds, if not hours or days.

Adenine
Pyrophosphate
Ribose
ATP

Entangled

Entanglement decay

Posner cluster
• Calcium

Entangled phosphates

① The biological molecule adenosine triphosphate (ATP) can release pyrophosphate, made from two phosphate molecules.

② Each phosphate carries a quantum spin, and the two phosphates can become entangled with each other.

③ Unprotected, the phosphate entanglement will decay, or decohere, in short order.

④ But if the phosphates are grouped together into protective clusters called Posner clusters, which are made of phosphate and calcium ions, the entanglement might survive for a longer time.

⑤ If a pair of entangled phosphates split into different Posner clusters, they will remain entangled even as the clusters transport them far from each other. In this way, the entanglement can be distributed over fairly long distances in the brain. This allows for the possibility of a quantum basis for brain function.

Lucy Reading-Ikkanda for Quanta Magazine

See also:
- https://www.quantamagazine.org/a-new-spin-on-the-quantum-brain-20161102/
- https://phys.org/news/2015-08-neural-qubits-quantum-cognition-based.html
- https://www.kitp.ucsb.edu/mpaf/quantum-brain

Universität Zürich UZH

37