

Road-mapping Biology

How the organization of biological knowledge
impacts revolution strategy

Adam Marblestone
Revolutionary Ventures 2015

Agenda

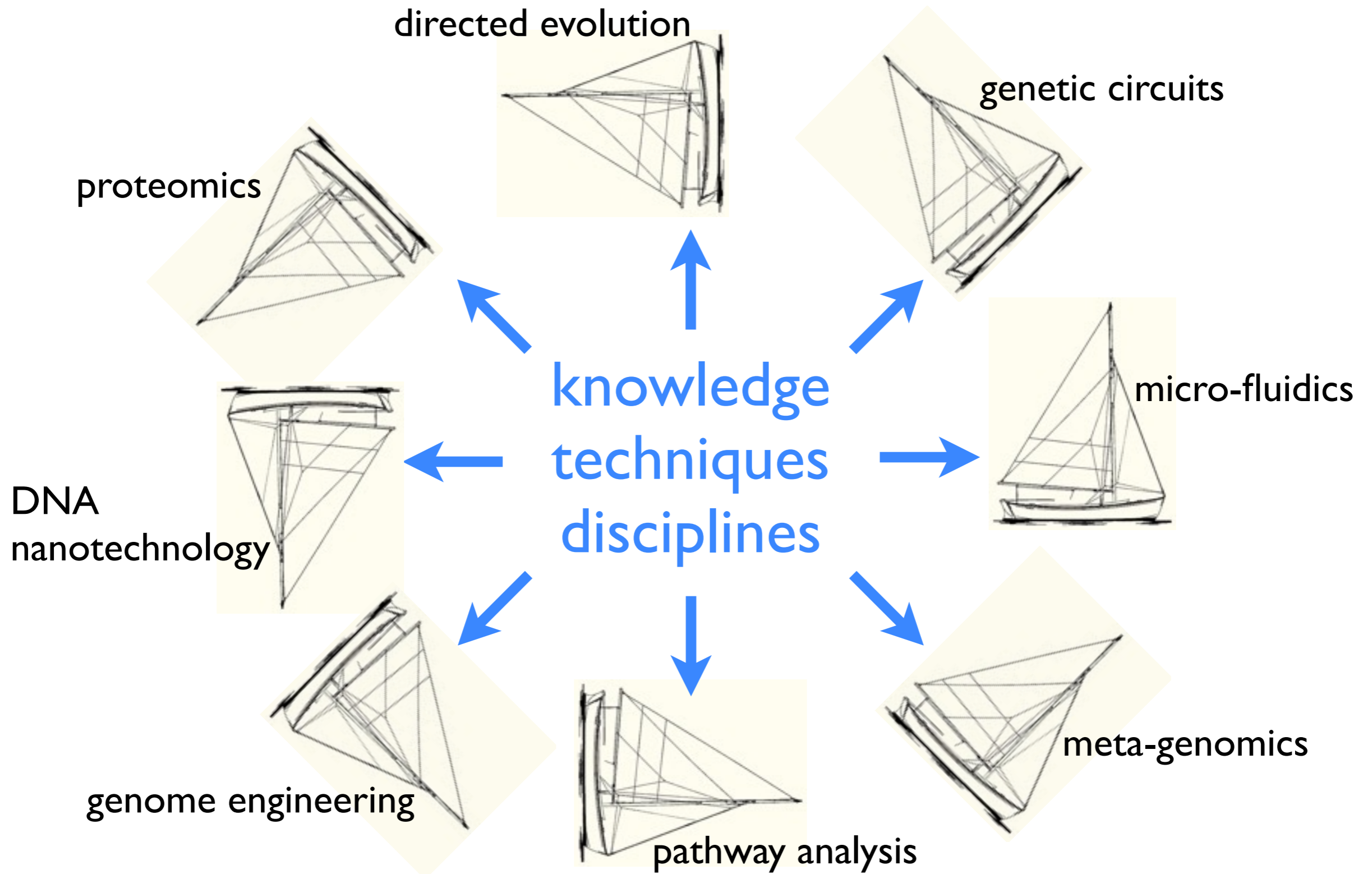
- What is a scientific roadmap?
- Structural inefficiencies in biology:
 - Why does biology need roadmaps?
 - Why are there so many hidden gems?
- Examples of roadmaps, and of hidden gems
- Improved software for mapping science

- **Roadmap:** A map of constraints, on the way towards a goal, and of potential workarounds for those constraints
- **Engineered Serendipity:** Biology breakthroughs depend on serendipity. We can make serendipity more likely by systematically surveying for hidden gems.

Why does biology need roadmaps?

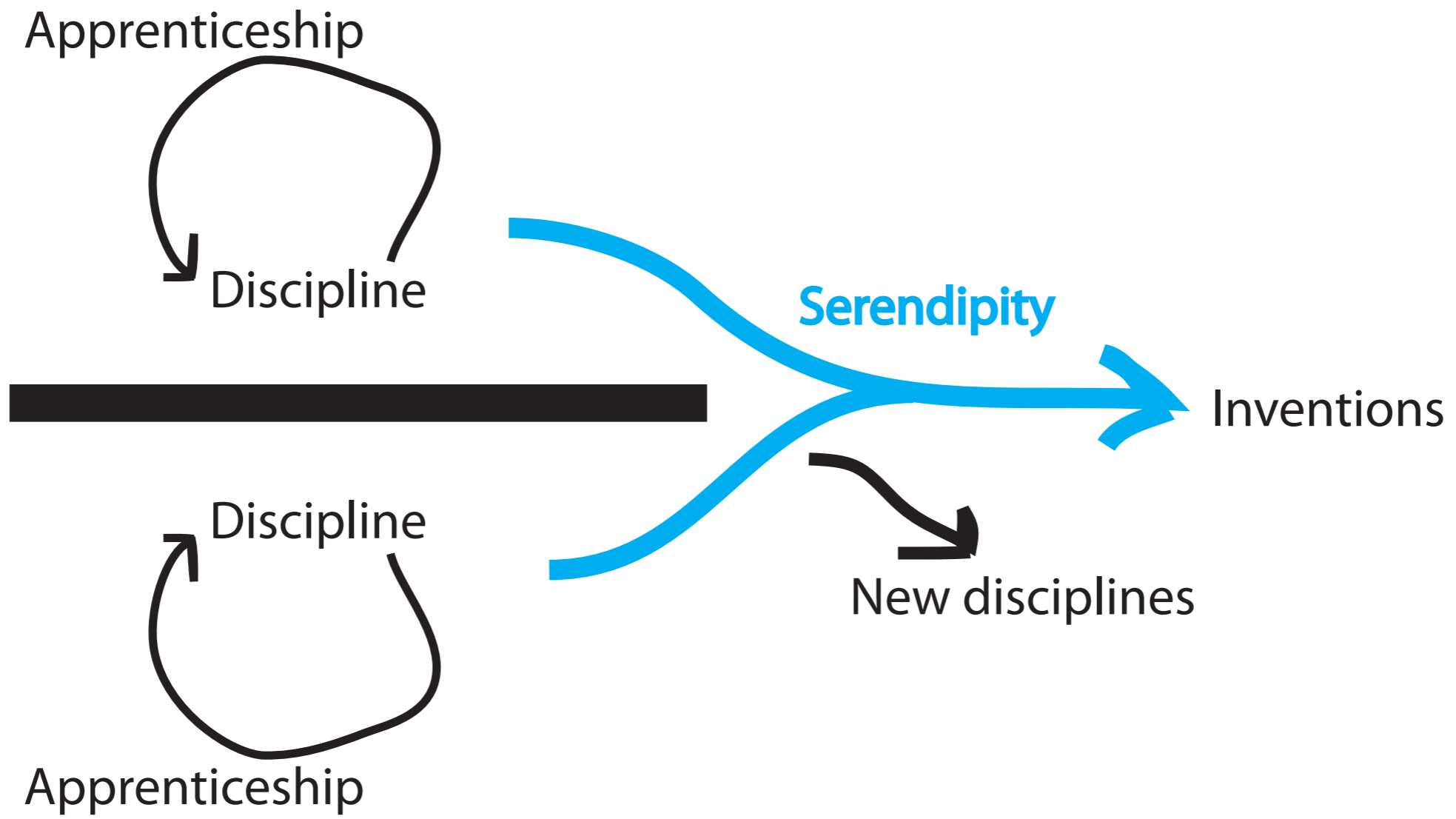
Why are there so many hidden gems?

“Science” does not have a plan for solving biotech grand challenges



A “sociological big bang”: does not scale well with problem complexity





The elephant in the room:

to make bio-technological quantum leaps
we must change how biology is done

**revolutions
by chance**



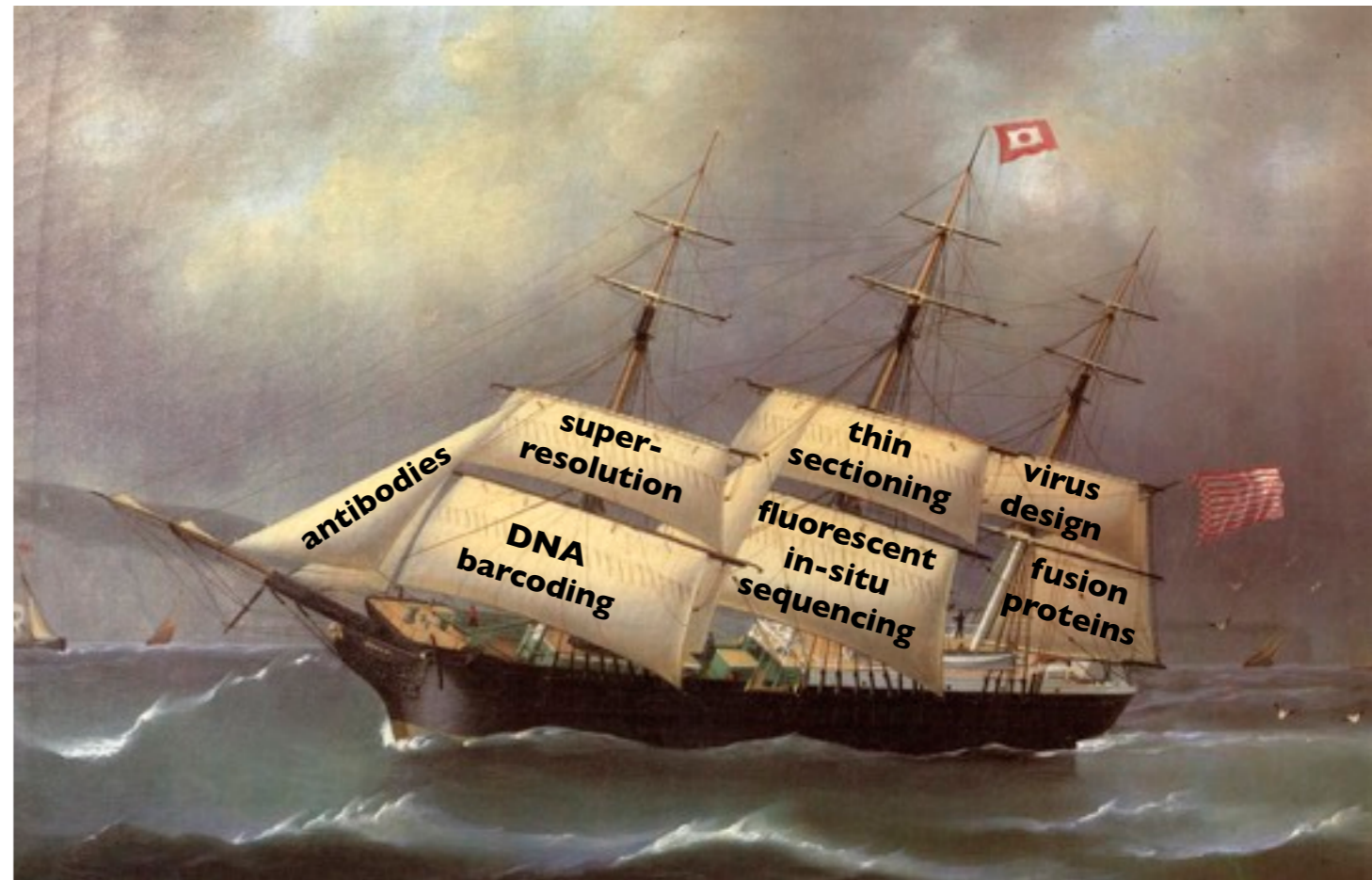
**revolutions
by design**

Two conspicuously missing elements in modern biology

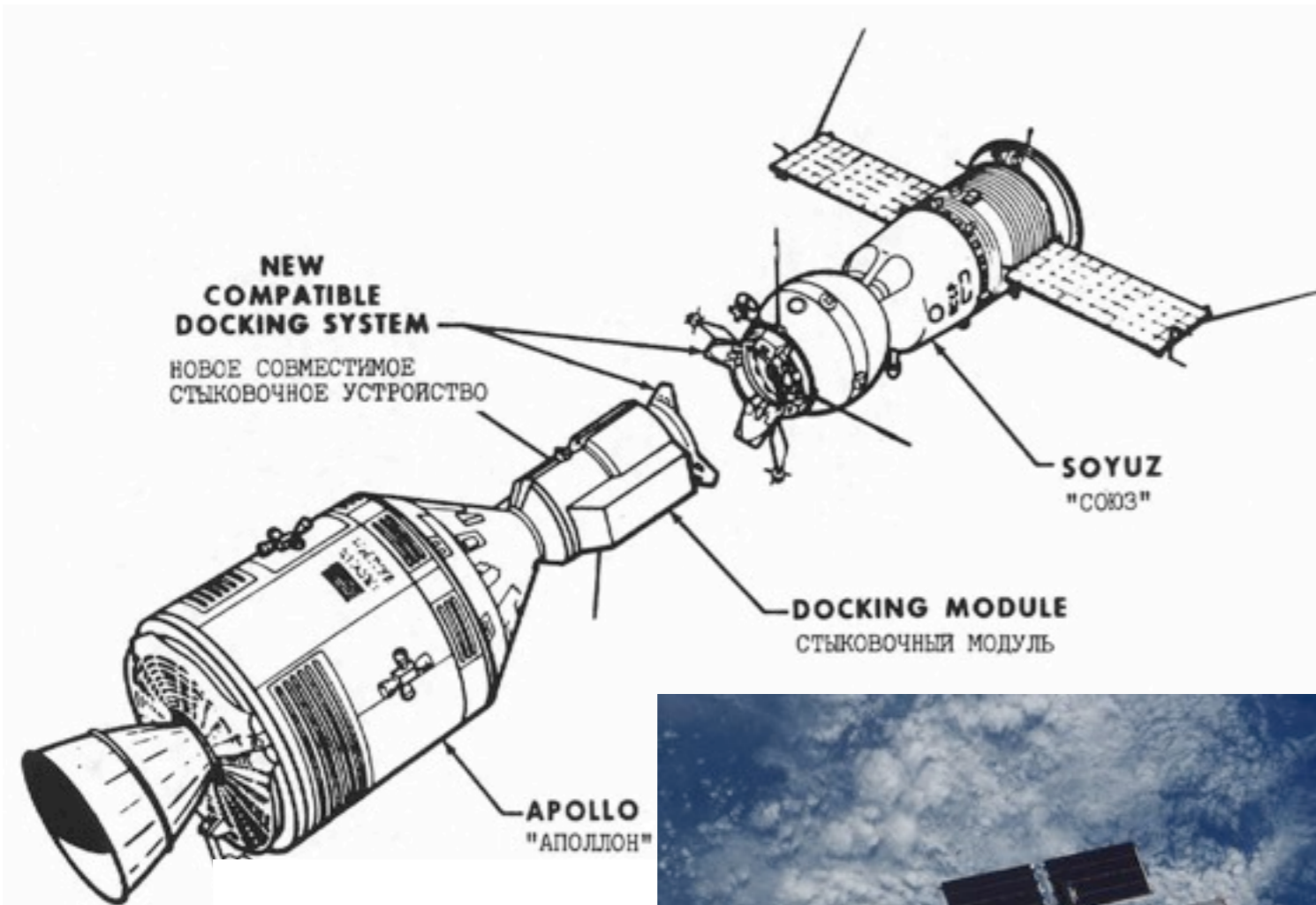
maps of entire domains



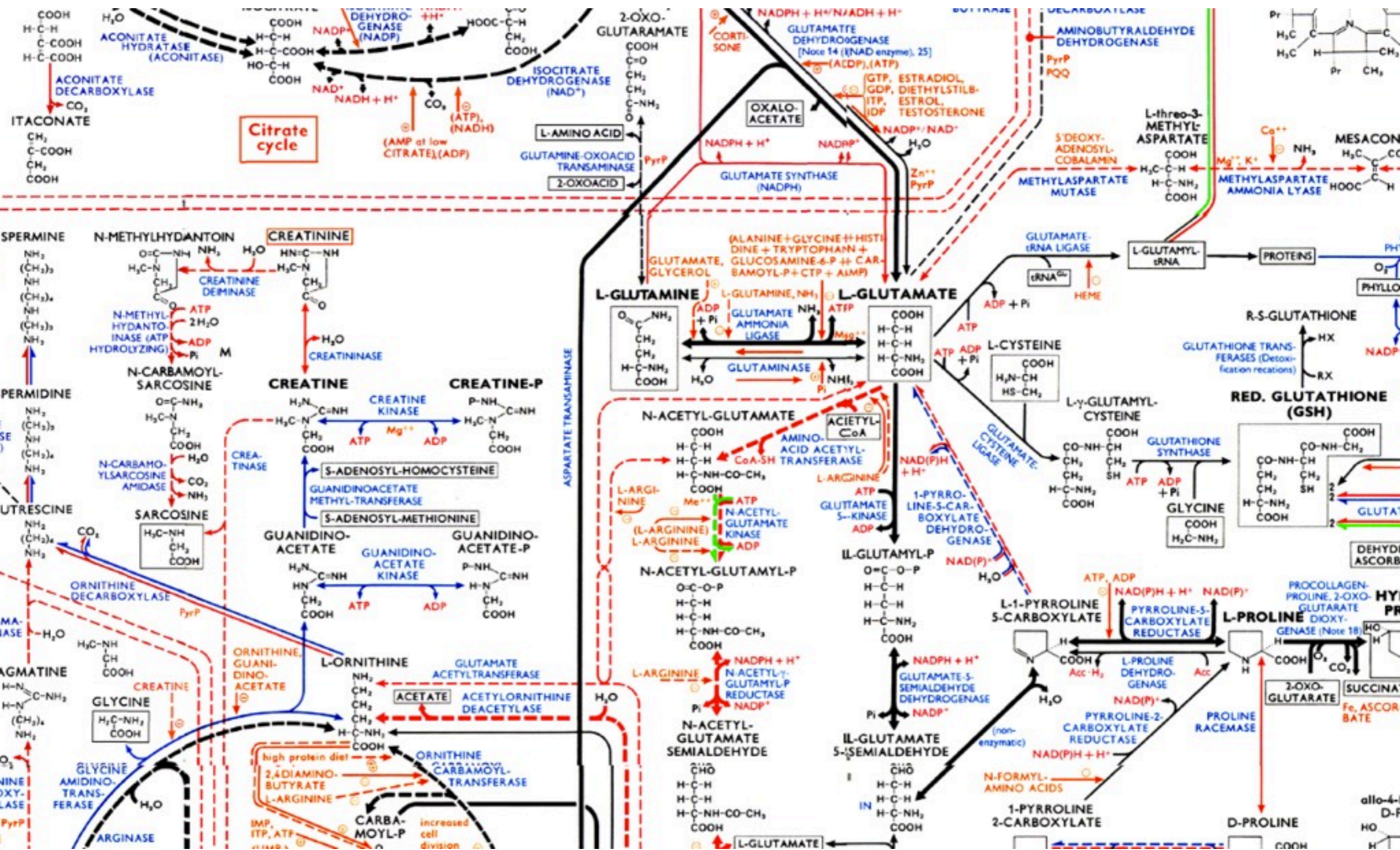
tech architectures / strategies



Well-known “systems engineering” methods in the sciences of *simplicity*...



How to extend to the sciences of complexity?



innovation “algorithm” & “business model” in bio is ripe for change...

How to create a bio-methods revolution

(only partially joking)

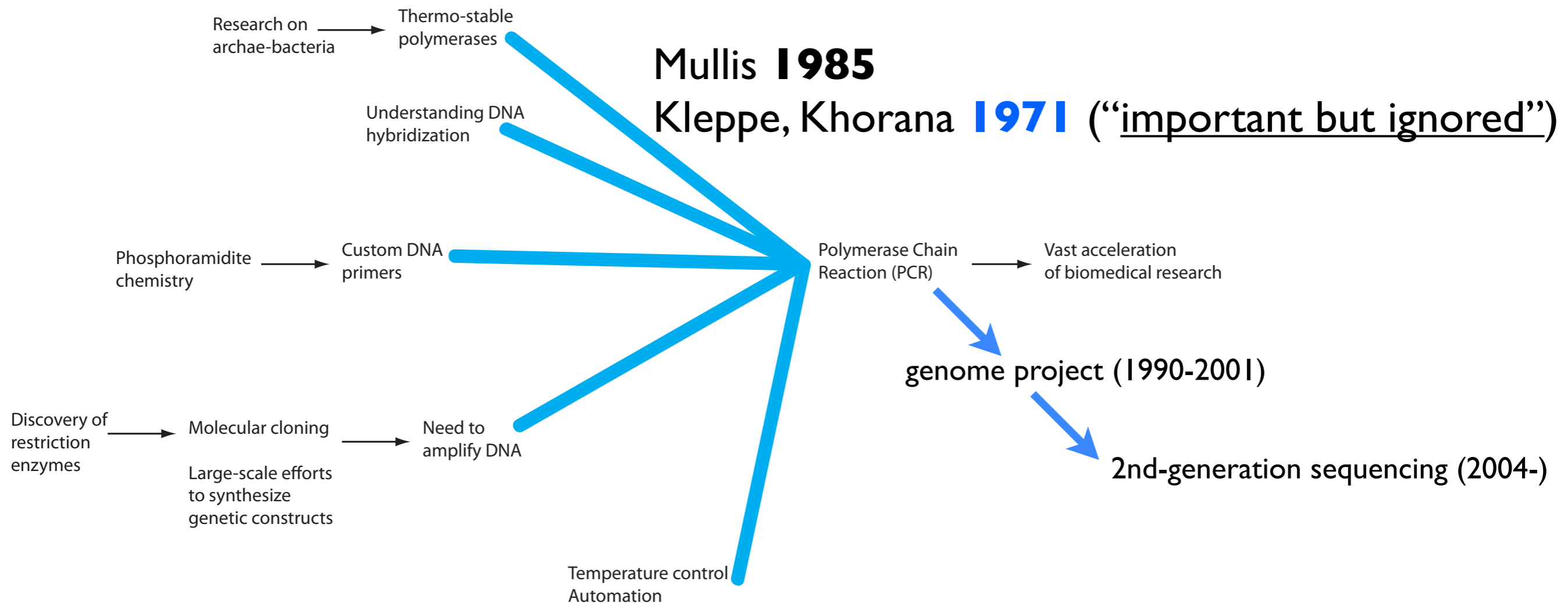
- 1) Understand landscape of *constraints* (physics, design robustness)
- 2) Find a *conceptual assumption-violating* workaround
- 3) Find a *hidden gem* (e.g., in nature) that implements this concept

#3 can be done via:

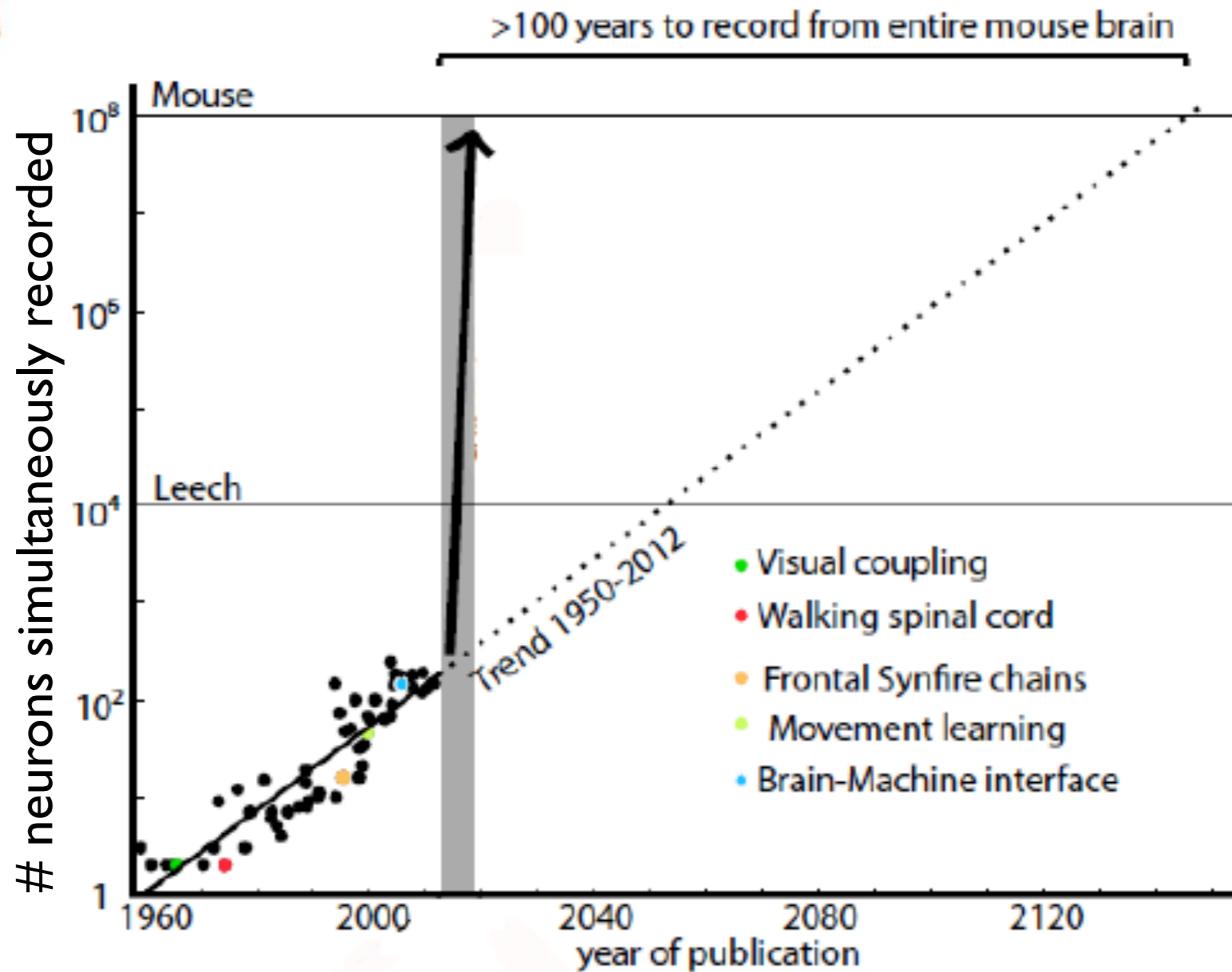
- a) *searching* existing knowledge and transplanting to new domain
- b) *screening* libraries of elements from nature (mining evolution)

PCR: a gem “hiding in plain sight”

Goal: make lots of copies of an arbitrary DNA sequence



Goal: record every “spike” from every neuron in a mammalian brain



nobody has written down a design that clearly
solves the problem
does not violate any laws of physics
does not severely damage the brain

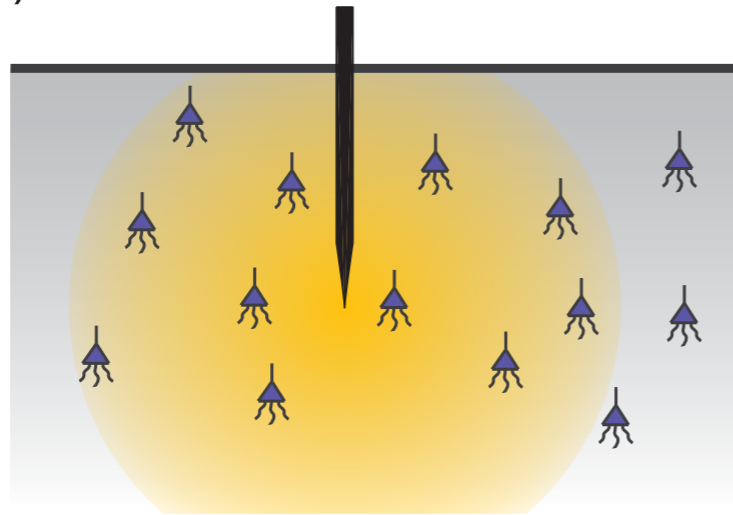
Landscape of constraints on brain activity mapping

< 2C temperature change:

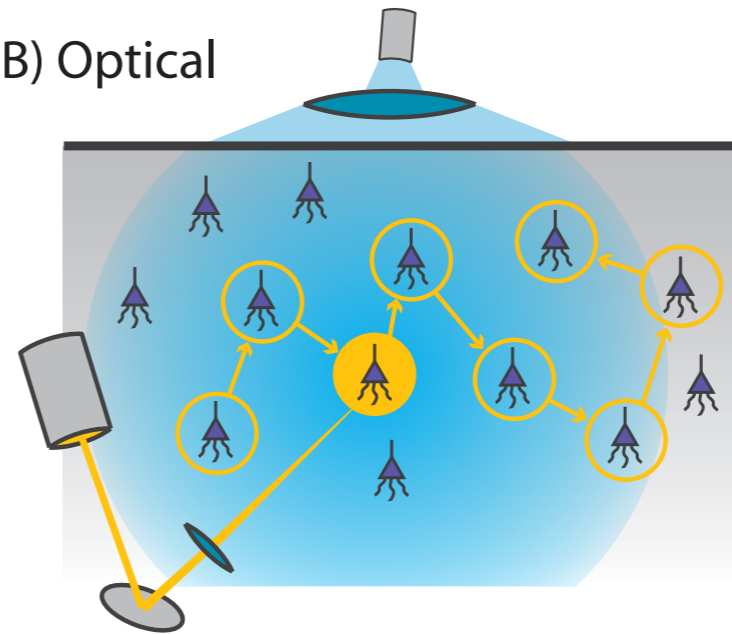
→ < **50 mW** steady-state power dissipation

< 1% tissue volume displacement

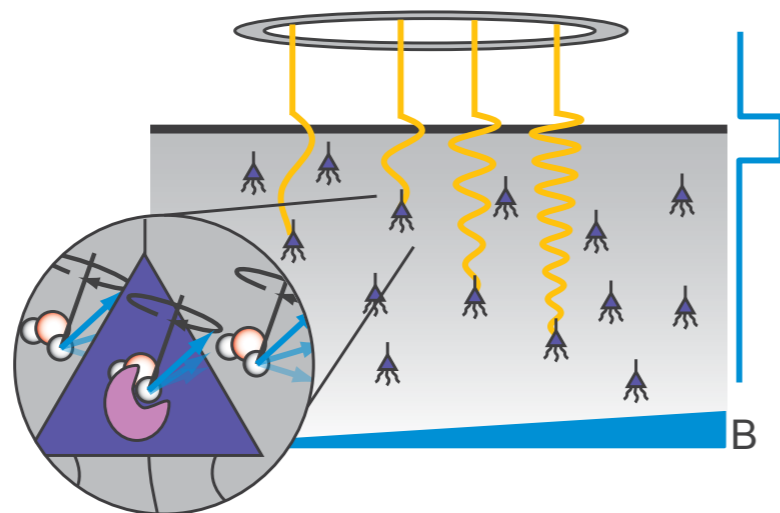
A) Electrical



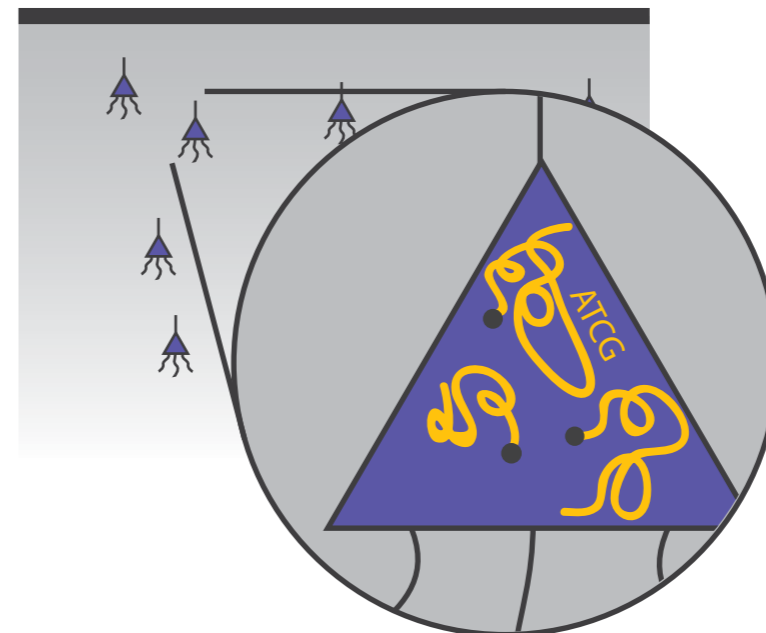
B) Optical



C) Magnetic Resonance



D) Molecular



Modality	Analysis Strategy	Assumptions	Conclusions
<i>Extracellular electrical recording</i>	<p>Compute minimal number of recorders based on max distance from recorder to recorded neuron</p> <p>Compute channel capacity limits to spike sorting</p>	<p>Decay profile of extracellular voltage</p> <p>Approximate noise levels at recording site</p>	<p>Maximum recording distance $r_{\max} \approx 100\text{--}200\mu\text{m}$ from electrode to neuron measured</p> <p>$\sim 10^5$ recording sites are required per mouse brain at current noise levels assuming perfect spike sorting</p> <p>$\sim 10^6$ recording sites are required at current noise levels at the physical limits of spike sorting</p> <p>$\sim 10^7$ recording sites are required using current spike sorting algorithms</p>
<i>Implanted electrical recorders</i>	<p>Compute power dissipation of electronic devices that digitally sample neuronal activity</p>	<p>Physical limit: $k_B T \ln(2)$/bit erased</p> <p>Practical limit: $\sim 10k_B T$/bit processed</p> <p>Current CMOS digital circuits: $> 10^5 k_B T$/bit processed</p>	<p>Requires 2–3 orders of magnitude increase in the power efficiency of electronics relative to current devices to scale to whole-brain simultaneous recordings</p> <p>Minimalist architectures could be developed to reduce local data processing overhead</p>
<i>Wireless data transmission</i>	<p>Compute tradeoff between power dissipation and channel bandwidth using information theory</p>	<p>Transmitter must supply enough power to overcome noise and path loss</p>	<p>Transmission at optical or near-optical frequencies is needed to achieve sufficient single-channel data rates using electromagnetic radiation. Radio-frequency (RF) electromagnetic transmission of whole-brain activity data draws excessive power due to bandwidth constraints</p> <p>Bandwidth cannot be split over multiple independent RF channels, but IR light or ultrasound may allow spatial multiplexing</p>
<i>Optical imaging</i>	<p>Relate the scattering and absorption lengths of optical wavelengths in brain tissue to signal-to-noise ratios for optical imaging</p>	<p>Approximate values of scattering and absorption lengths as a function of wavelength</p>	<p>Light scattering imposes severe constraints, but strategies exist which could negate the effects of scattering, such as implantable optics, infrared indicators, signal modulation, and online inversion of the scattering matrix</p>
<i>Multi-photon optical imaging</i>	<p>Compute minimum total excitation light power to excite multi-photon transitions from indicators within each neuron in every imaging frame</p>	<p>Approximate values of multi-photon cross-sections</p> <p>Pulse durations similar to those currently used in multi-photon imaging</p>	<p>Whole-brain multi-photon excitation will over-heat the brain except in very short experiments, unless ultra-high-cross-section indicators are used</p>
<i>Beam scanning microscopies</i>	<p>Calculate device and indicator parameters necessary for fast beam repositioning and signal detection</p>	<p>Fast optical phase modulators could reposition beams at ~ 1 GHz switching rates</p> <p>Fluorescence lifetimes in the 0.1–1.0 ns range</p>	<p>Beam repositioning time limits the speed of current systems but these are far from the physical limits</p> <p>Fluorescence lifetimes of indicators constrain design of ultra-fast scanning microscopies</p>
<i>Magnetic resonance imaging</i>	<p>Calculate spatial and temporal resolution of MRI based on spin relaxation times and spin diffusion</p>	<p>Proton MRI using tissue water</p> <p>Approximate T_1 and T_2 relaxation times and self-diffusion times for tissue water</p>	<p>Proton MRI is limited by the T_1 relaxation time of water to ~ 100 ms temporal resolution and by the self-diffusion of water to spatial resolutions of $\sim 40\mu\text{m}$. T_1 pre-mapping could allow T_2 contrast on a ~ 10 ms timescale. Achieving these limits for functional imaging requires going beyond BOLD contrast</p>
<i>Ultrasound</i>	<p>Calculate spatial resolution, signal strength and bandwidth limits on ultrasound imaging</p>	<p>Speed of sound in brain</p> <p>Attenuation length of ultrasound in brain</p>	<p>Attenuation of ultrasound by brain tissue and bone may be prohibitive at the ~ 100 MHz frequencies needed for single-cell resolution ultrasound imaging</p> <p>Ultrasound may be viable for spatially multiplexed data transmission from embedded devices [70]</p>
<i>Molecular recording</i>	<p>Compute metabolic load and volume constraint for rapid synthesis of large nucleic acid polymers</p> <p>Evaluate temporal resolution in simulated experiments using kinetic models [6]</p>	<p>Polymerase biochemical parameter ranges</p> <p>Metabolic requirements of genome replication</p>	<p>Molecular recording devices appear to fall within physical limits but their development poses multiple major challenges in synthetic biology</p> <p>Synchronization or time-stamping mechanisms are required for temporal resolution to approach the millisecond scale</p>

many electrodes are needed

embedded electronics are too power-hungry, at present

use IR or ultrasound, not RF, for data-transmission

light scattering can potentially be overcome in several ways

multi-photon optics is too dissipative

requires many parallel scanned beams

MRI needs new contrast mechanisms

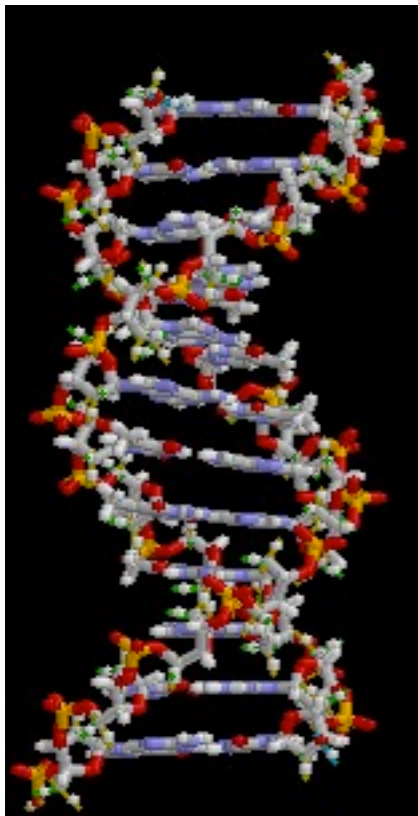
ultrasound is potentially powerful

molecular recording is possible but hard

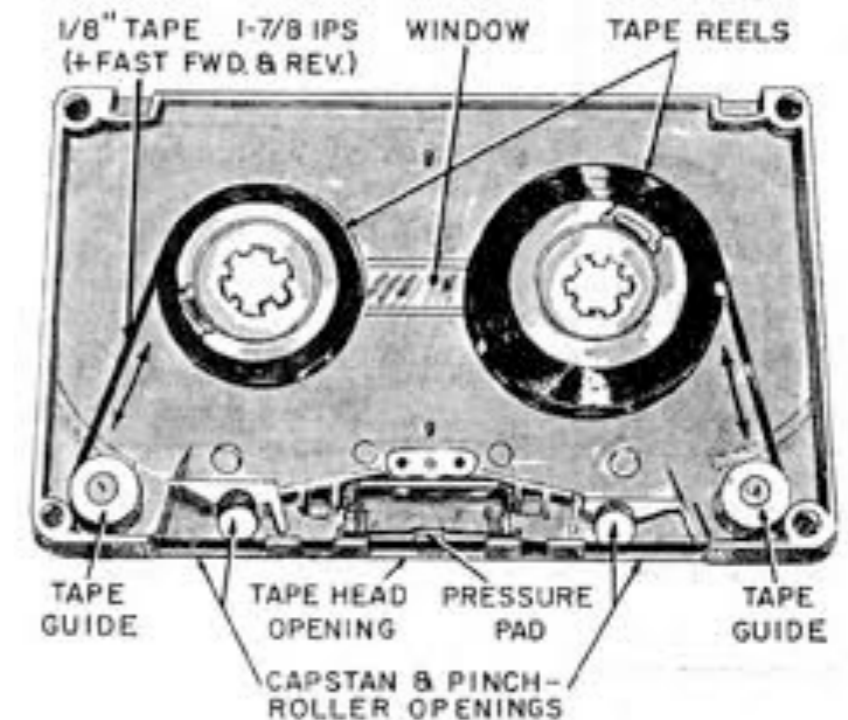
Need for a conceptual work-around



Assumption-violating concept: ***what if each cell could record its own activity?***



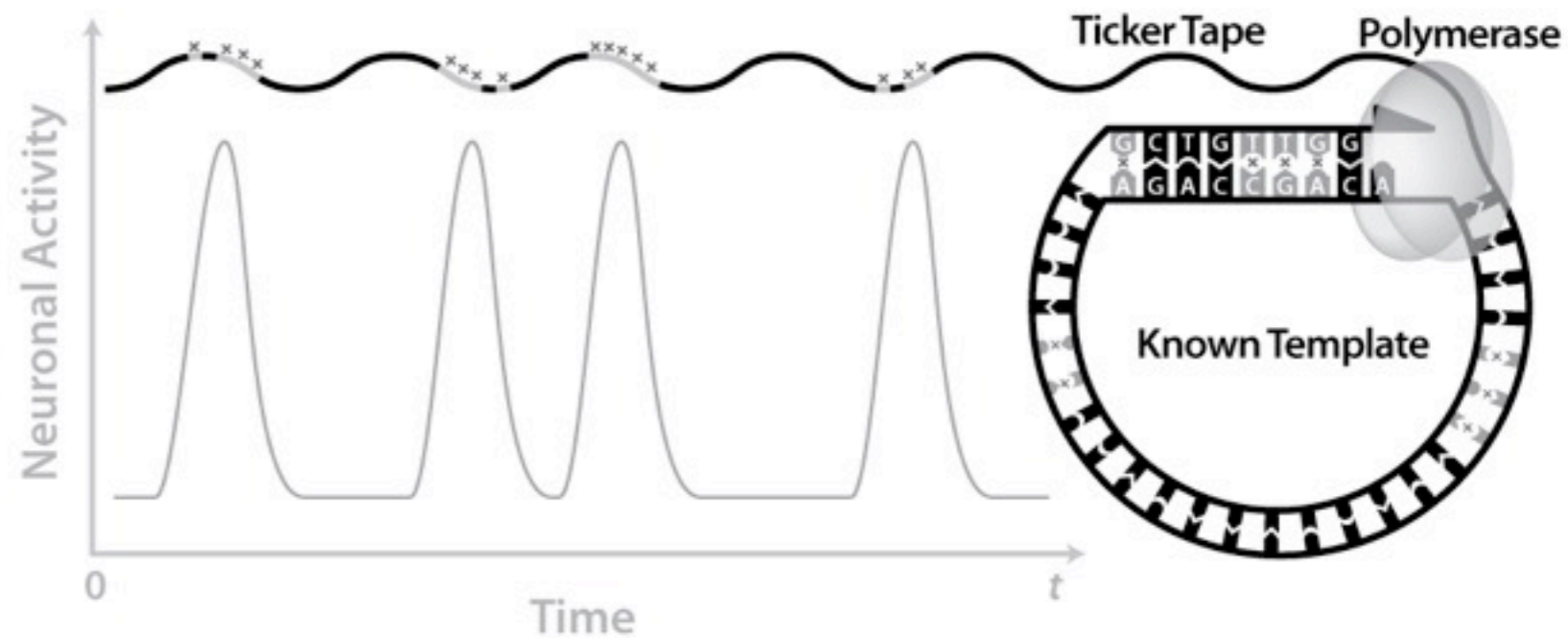
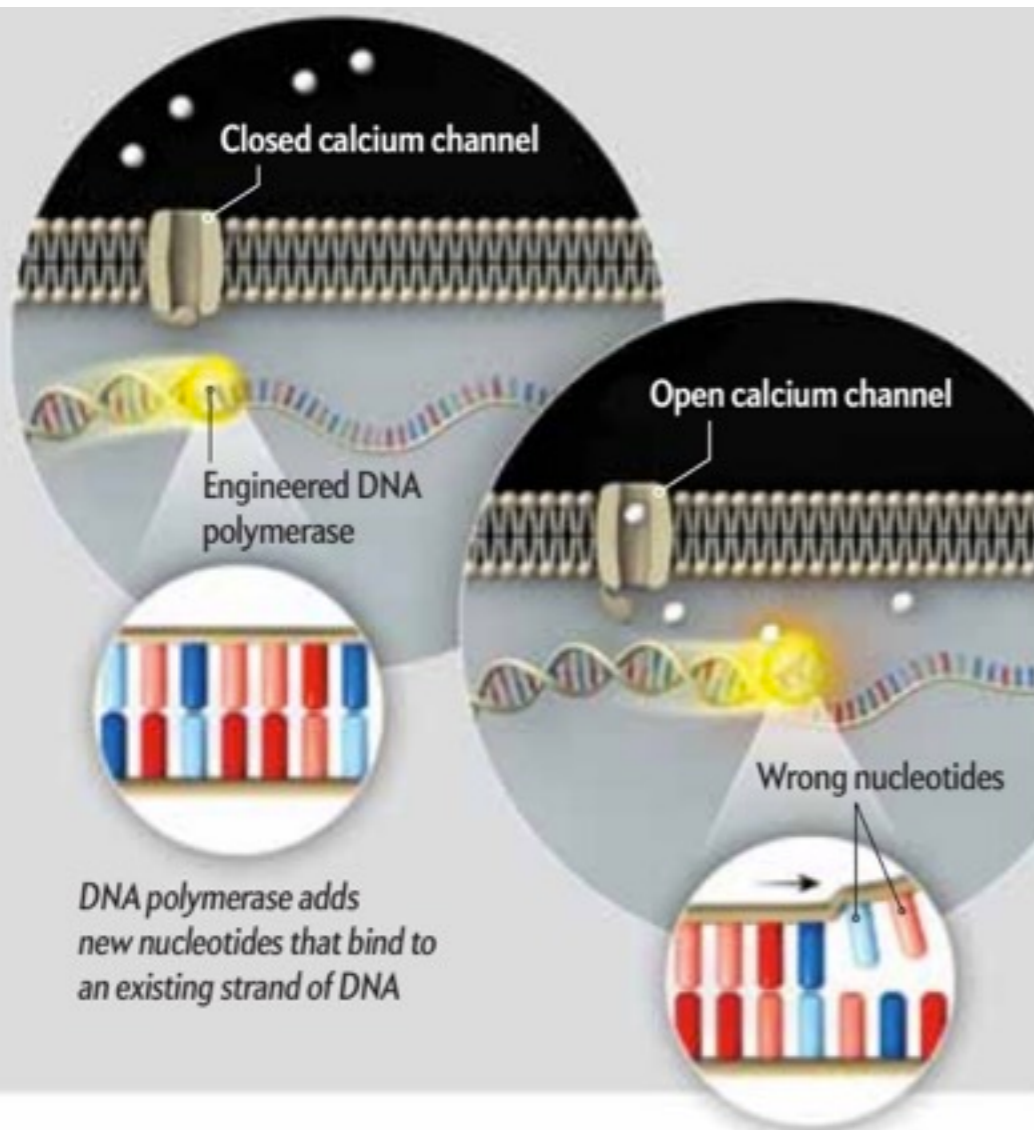
=
0101001010101010010101
0101010101001010110111
0010110101010111010011
=



\$9M NIH grant to pursue this

Molecular implementation of the concept: *a molecular recording device in each neuron*

Encode information via control of DNA polymerase copying **error rate**



Mapping out the space of *theories* about the brain

Computation	Algorithmic/ representational realization	Neural implementation(s)	Brain location(s)
Rapid perceptual classification	Receptive fields, pooling and local contrast normalization ^{51,55}	Hierarchies of simple and complex cells ⁵⁶	Visual system
Complex spatiotemporal pattern recognition	Bayesian belief propagation ^{19,57}	Feedforward and feedback pathways in cortical hierarchy ¹⁹	Sensory hierarchies
Learning efficient coding of inputs	Sparse coding ⁵⁸	Thresholding and local competition ⁵⁹	Sensory and other systems
Working memory	Continuous or discrete attractor states in networks ^{60,61}	Persistent activity in recurrent networks ⁶²	Prefrontal cortex
Decision making	Temporal-difference reinforcement learning algorithms ^{63,64} ; actor-critic models ⁶⁵	Cortically implemented Bayesian inference networks combined with temporal difference reinforcement learning via the dopamine system and action selection systems in the basal ganglia ⁶⁶	Prefrontal cortex
	Winner-take-all networks ⁶⁷	Recurrent networks coupled via lateral inhibition ⁶⁷	Prefrontal cortex
Gating of information flow	Context-dependent tuning of activity in recurrent network dynamics ⁶⁸	Recurrent neural networks implementing line attractors and selection vectors ⁶⁸	Prefrontal cortex
	Shifter circuits ⁶⁹	Divergent excitatory relays and input-selective shunting inhibition in dendrites ⁶⁹	Visual system
Gain control	Divisive normalization ⁵²	Shunting inhibition in networks or balanced background synaptic excitation and inhibition ⁷⁰	Common across many cortical areas
Sequencing of events over time ⁷¹	Feed-forward cascades; Serial working memories ⁷²	Synfire chains ⁷³⁻⁷⁵ ; Thalamo-cortico-striatal loops ^{76,77}	Common across many cortical areas
Representation and transformation of variables	Population coding ⁷⁸	Time-varying firing rates of cosine-tuned neurons representing dot products with encoding vectors	Motor cortex
Variable binding	Holographic reduced representations ^{49,79}	Circular convolution of vectors represented by neural population codes	Cortical areas involved in sequential or symbolic processing
	Dynamic binding ^{80,81}	Neural synchronization ⁸²	

with Gary Marcus
(NYU + Allen Institute)
and Tom Dean (Google)

Can software tools accelerate the uptake of cross-disciplinary knowledge, helping us find the hidden gems?

Automatically learning the “meanings” of science words

Word2Vec model trained on 150k PubMed abstracts

```
prompt> dopamin
dopamine          dopaminergic      dopaminergically dopaminoceptive dopaminomimetic dopaminomimetics

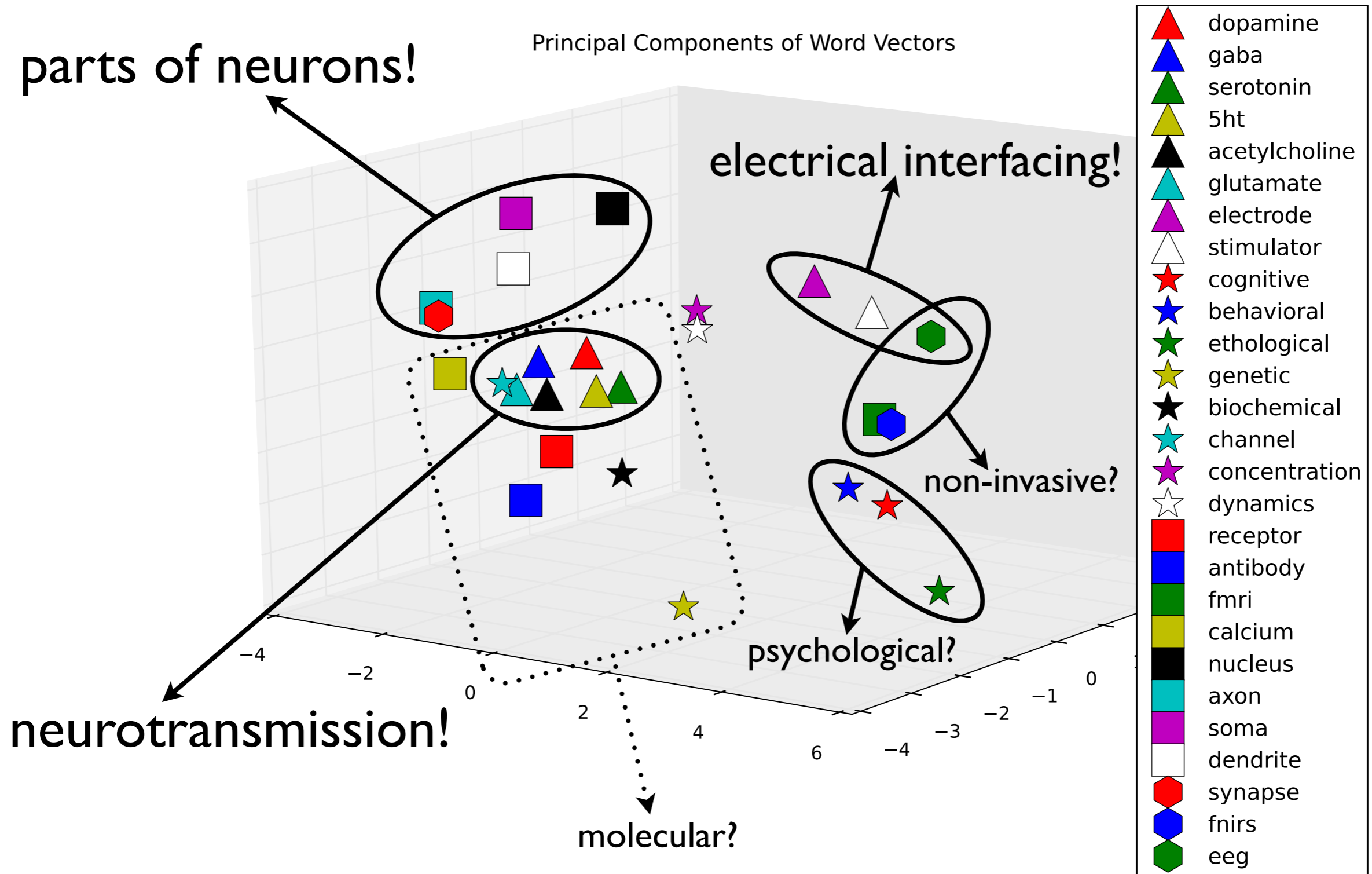
prompt> dopamine

Similar words:
-----
serotonergic
nsda
5ht
dopamine
noradrenergic
monoamines
nigro
striatal
dopaminergic
cholinergic
monoamine
hydroxytryptamine
nigrostriatal
reuptake
serotonin
nigral
norepinephrine
midbrain
msns
sert
striatum
noradrenalin
noradrenaline
daergic
mamph
```

statistical model was only fed raw text, and knows nothing about neuroscience (or anything), yet it “discovers” neurotransmitters

Automatically learning the “meanings” of science words

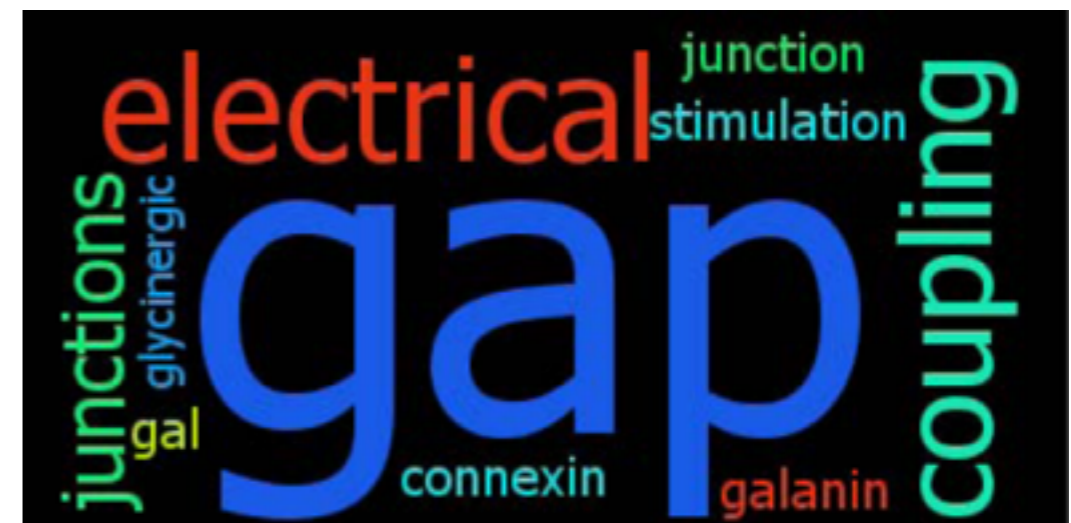
Word2Vec model trained on 150k PubMed abstracts



Identifying the sub-fields/sub-topics:
unsupervised document classification model



gaba is an inhibitory neurotransmitter

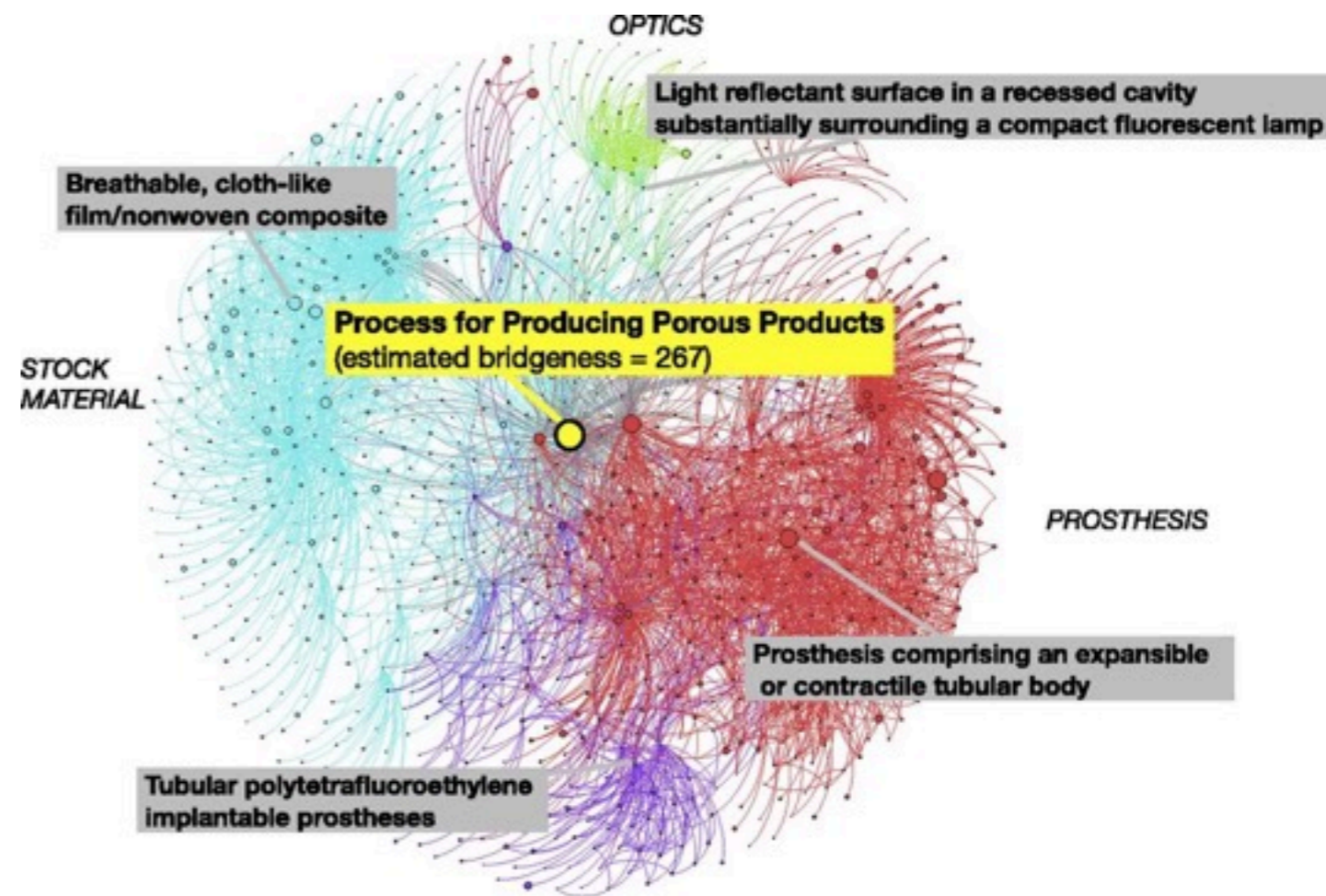
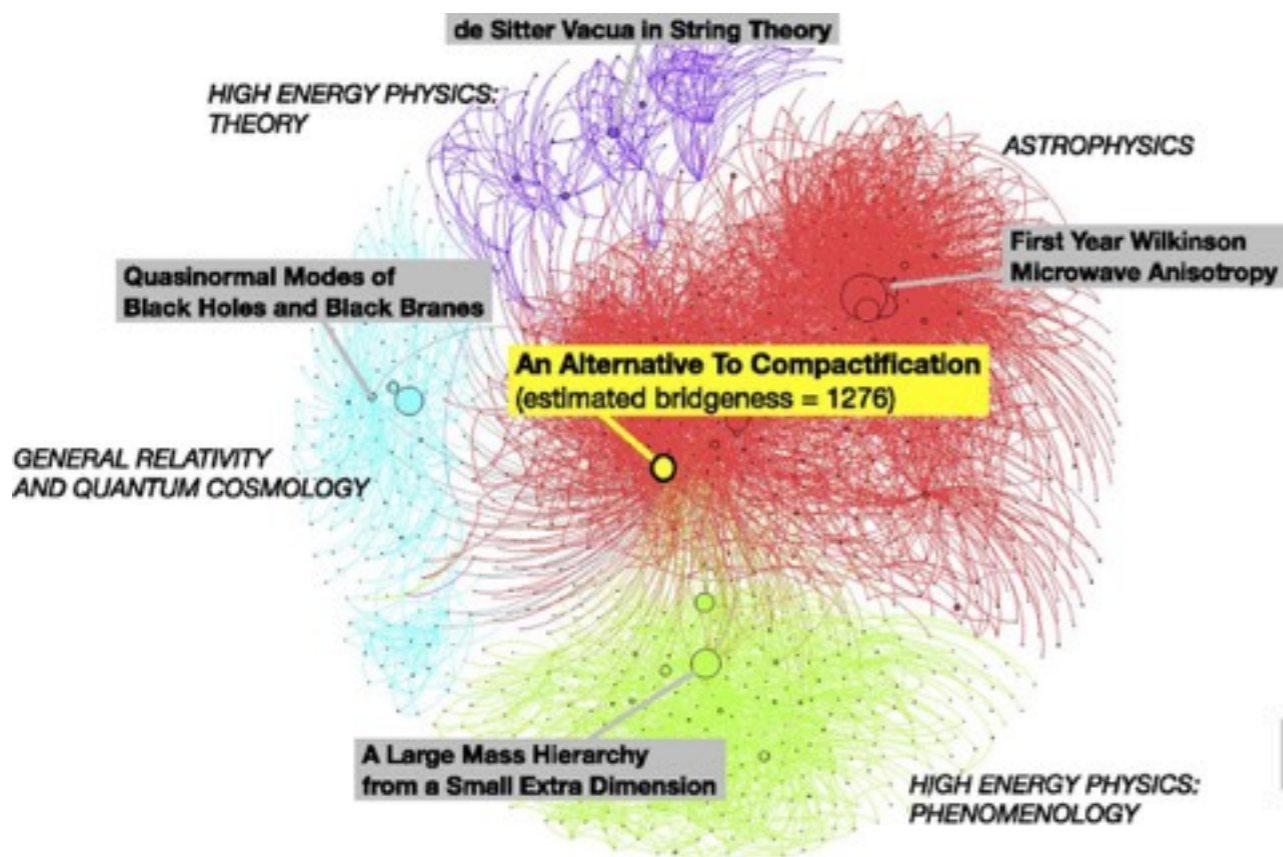


connexin is a gap-junction protein



the *dopamine* system lives in the *substantia nigra* region of the *striatum*

Identifying the sub-fields/sub-topics: community-detection on the citation graph



Efficient discovery of overlapping communities in massive networks


Prem K. Gopalan¹ and David M. Blei

Department of Computer Science, Princeton University, Princeton, NJ 08540

Identifying the sub-fields/sub-topics: community-detection on the citation graph

The first large-scale publicly available citation graph

Microsoft Academic Graph



The Microsoft Academic Graph is a heterogeneous graph containing scientific publication records, citation relationships between those publications, as well as authors, institutions, journals and conference "venues" and fields of study. This data is available as a set of zipped text files stored in Microsoft Azure blob storage and available via HTTP.

Update: August 31 2015

An updated version of the Microsoft Academic Graph has been released. This version includes more papers and citations, improved conflation of duplicate entities, and is provided in smaller downloadable chunks to improve access.

To download the data you need to first agree to the [terms of use](#).

☐ I agree to abide by the [terms of use](#) for the Microsoft Academic Graph. [Get the data!](#)

Towards human-computer synergy for accelerated science

No. 4356 April 25, 1953 NATURE 737

equipment, and to Dr. G. E. B. Deacon and the captain and officers of R.R.S. *Discovery II* for their part in making the observations.

¹Young, P. E., Gerani, H., and Jevons, W., *Phil. Mag.*, **40**, 149 (1925).

²Longuet-Higgins, M. S., *Mon. Not. Roy. Astro. Soc., Cambro. Symp.*, **3**, 265 (1950).

³Von ARK, W. S., *Weeks Note Papers in Phys. Ozeanog. Meteor.*, **11** (3) (1950).

⁴Ekman, V. W., *Arkiv. Mat. Astron. Fysik. (Stockholm)*, **2** (11) (1950).

MOLECULAR STRUCTURE OF NUCLEIC ACIDS

A Structure for Deoxyribose Nucleic Acid

WE wish to suggest a structure for the salt of deoxyribose nucleic acid (D.N.A.). This structure has novel features which are of considerable biological interest.

A structure for nucleic acid has already been proposed by Pauling and Corey¹. They kindly made their manuscript available to us in advance of publication. Their model consists of three intertwined chains, with the phosphates near the fibre axis, and the bases on the outside. In our opinion, this structure is unsatisfactory for two reasons: (1) We believe that the material which gives the X-ray diagrams is the salt, not the free acid. Without the acidic hydrogen atoms it is not clear what forces would hold the structure together, especially as the negatively charged phosphates near the axis will repel each other. (2) Some of the van der Waals distances appear to be too small.

Another three-chain structure has also been suggested by Fraser (in the press). In his model the phosphates are on the outside and the bases on the inside, linked together by hydrogen bonds. This structure as described is rather ill-defined, and for this reason we shall not comment on it.

We wish to put forward a radically different structure for the salt of deoxyribose nucleic acid. This structure has two helical chains each coiled round the same axis (see diagram). We have made the usual chemical assumptions, namely, that each chain consists of phosphate di-ester groups joining β -D-deoxy-ribofuranose residues with 3',5' linkages. The two chains (but not their bases) are related by a dyad perpendicular to the fibre axis. Both chains follow right-handed helices, but owing to the dyad the sequences of the atoms in the two chains run in opposite directions. Each chain loosely resembles Purrberg's² model No. 1; that is, the bases are on the inside of the helix and the phosphates on

is a residue on each chain every 3.4 Å. in the z-direction. We have assumed an angle of 36° between adjacent residues in the same chain, so that the structure repeats after 10 residues on each chain, that is, after 34 Å. The distance of a phosphorus atom from the fibre axis is 10 Å. As the phosphates are on the outside, cations have easy access to them.

The structure is an open one, and its water content is rather high. At lower water contents we would expect the bases to tilt so that the structure could become more compact.

The novel feature of the structure is the manner in which the two chains are held together by the purine and pyrimidine bases. The planes of the bases are perpendicular to the fibre axis. They are joined together in pairs, a single base from one chain being hydrogen-bonded to a single base from the other chain, so that the two lie side by side with identical z-co-ordinates. One of the pair must be a purine and the other a pyrimidine for bonding to occur. The hydrogen bonds are made as follows: purine position 1 to pyrimidine position 1; purine position 6 to pyrimidine position 6.

If it is assumed that the bases only occur in the structure in the most plausible tautomeric forms (that is, with the keto rather than the enol configurations) it is found that only specific pairs of bases can bond together. These pairs are: adenine (purine) with thymine (pyrimidine), and guanine (purine) with cytosine (pyrimidine).

In other words, if an adenine forms one member of a pair, on either chain, then on these assumptions the other member must be thymine; similarly for guanine and cytosine. The sequence of bases on a single chain does not appear to be restricted in any way. However, if only specific pairs of bases can be formed, it follows that if the sequence of bases on one chain is given, then the sequence on the other chain is automatically determined.


It has been found experimentally^{3,4} that the ratio of the amounts of adenine to thymine, and the ratio of guanine to cytosine, are always very close to unity for deoxyribose nucleic acid.

It is probably impossible to build this structure with a ribose sugar in place of the deoxyribose, as the extra oxygen atom would make too close a van der Waals contact.

The previously published X-ray data^{3,4} on deoxyribose nucleic acid are insufficient for a rigorous test of our structure. So far as we can tell, it is roughly compatible with the experimental data, but it must be regarded as unproved until it has been checked against more exact results. Some of these are given in the following communications. We were not aware of the details of the results presented there when we devised our structure, which rests mainly though not entirely on published experimental data and stereochemical arguments.

It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material.

Full details of the structure, including the conditions assumed in building it, together with a set of co-ordinates for the atoms, will be published elsewhere.



Icons: Note, Highlight, Screenshot, Flag

Publication

A Structure for Deoxyribose Nucleic Acid

Francis Crick James Watson

234 10634 Tools

Conversations

Re: A Structure for Deoxyribose Nucleic Acid updated 20 minutes ago 12k

Re: Physical Principles for Scalable Neural Rec... updated 9 years ago 234

Re: Millisecond-timescale, genetically targeted ... updated 9 years ago 1

Re: A Structure for Deoxyribose Nucleic Acid updated 20 minutes ago 12k

Re: Physical Principles for Scalable Neural Rec... updated 9 years ago 234

Re: Millisecond-timescale, genetically targeted ... updated 9 years ago 1

With Juan Batiz-Benet, Richard Littauer, Ed Boyden

Towards human-computer synergy for accelerated science

Re: A Structure for Deoxyribose Nucleic Acid

B

I

This paper is really good because it launched a new field after coming up with the [greatest discovery](#) of all time. This paper is really good because it launched a [new field](#) after coming up with the greatest discovery of all time. This paper is really good. Here is a web link, here is a link [into the paper](#), here is a link to a [@Person](#), here is a link to a [Conversation](#), here is a link to [a figure](#), and i think this paper is [#interesting](#).

+ invite others to the conversation

Read

Juan Benet
juan@benet.ai

Initiator

Start new conversation

Rosalind Franklin

+ Add

Share

Juan Benet
juan@benet.ai

Read can view conversation Initiator

Francis Crick
forick@mro.ao.uk

Write can add content to the conversation Participate

James Watson
jwat@mro.ao.uk

Share can invite others to the conversation

Read