

## Neuroscience Keeps Solving the Same Problems Twice

*A structured claim-dependency layer for science, and an experiment to test whether it accelerates convergence.*

**Introduction** - Here is a story about how a scientific field can lose two decades to a question it had the tools to settle much sooner.

In 1986, Apostolos Georgopoulos and colleagues published one of the most cited results in motor neuroscience [1]. They recorded from neurons in monkey motor cortex during reaching and found that each neuron fired most for one preferred direction. Weight each neuron's preferred direction by its firing rate, sum the result, and you get a "population vector" that predicts the direction of the upcoming reach. It was elegant. It felt explanatory. Motor cortex represents direction.

Fourteen years later, Emanuel Todorov showed this inference was almost certainly wrong [2]. Not the data. The interpretation. Motor cortex sends signals to muscles. Muscles move a limb with particular lengths, masses, and joint configurations. Todorov built a model of this biomechanical chain and showed that commands to muscles will naturally produce neural activity correlated with movement direction, simply because of arm geometry. Sanger had shown earlier, on purely mathematical grounds, that a population vector can always be found under very general assumptions [3]. Direction tuning was not a neural code for direction. It was a byproduct of controlling a physical limb.

You might expect this to have forced a course correction. It didn't. The debate about what motor cortex neurons "encode" continued for another twenty years. When the field eventually shifted toward dynamical systems models in the 2010s, many of the same confusions reappeared in new vocabulary. Todorov's deeper question, what you can actually infer from neural correlations once you account for biomechanics, was never cleanly resolved.

I run into this pattern constantly. My own work develops methods for aligning neural representations across subjects and species [4], which means I regularly collaborate with experimentalists in motor cortex, songbird vocal production, and prefrontal decision-making. Before I can compare results, I spend weeks doing something that feels absurd for a mature science. I reconstruct, from scratch, which published claims are genuinely in tension, which ones merely appear contradictory because they depend on different preprocessing choices, and which reflect real disagreements about the underlying biology. This is slow even for specialists. For newcomers it is nearly impossible.

The standard explanation is incentives. Scientists are rewarded for novelty, not synthesis. That matters, but it's incomplete. Consider a different explanation. Science has become very good at distributing papers and very bad at storing structured knowledge about what's inside those papers. The bottleneck is not just motivation. It's infrastructure.

Today's scientific databases can tell you that Paper A cites Paper B and that both concern motor cortex. What they cannot tell you is whether Paper A's central inference depends on a particular smoothing kernel, whether Paper C later showed that kernel changes the result qualitatively, or whether Paper D's apparent disagreement with Paper A is about theory, preprocessing, or task design. When the only durable unit of scientific knowledge is the paper, every kind of relationship gets flattened into the same object. Support, challenge, scope restriction, and methodological dependence all become a citation. That is an extraordinarily poor basis for collective memory.

**The recycling pattern** - Motor cortex is a useful case study because its debates are well documented and they keep restarting. Force in the 1960s. Direction in the 1980s. Kinematics versus kinetics. Then representation versus dynamics. Each generation partly absorbed the previous one, but rarely made the inheritance explicit.

Why? Theoretical critiques fail to propagate. Todorov and others showed that once you ask not "what variable do neurons encode?" but "what control law does the circuit implement, given sensory feedback and the biomechanics of the body?" [5], many of the older disputes collapse. But the older framing persisted because it was operationally easier. Computing a population vector is straightforward. Building a feedback control model is hard. The easier analysis generates results faster, even if it answers a less well-posed question.

You can watch a version of this happening right now, in real time. In 2012, Churchland et al. [6] showed that reaching-related neural activity is well described by rotational dynamics, a result that reshaped the field. Since then, a series of papers have each changed what that result means, but without anyone connecting them into a coherent picture. Sauerbrei et al. [7] showed cortical activity depends on continuous thalamic input, challenging the idea that the dynamics are autonomous. Kalidindi et al. [8] showed rotational structure is what you'd expect

from a feedback controller, not necessarily evidence for a central pattern generator. Elsayed and Cunningham [9] showed that temporal autocorrelations preserved by standard smoothing can produce spurious rotational structure, making the detection of rotations sensitive to preprocessing choices. And Suresh et al. [10] showed the dynamics seen during reaching do not appear during grasping. Each finding should change how you interpret the 2012 result. But there is no standard resource connecting them, so the field carries on as if each paper exists in isolation.

**Why existing tools are not enough** - You might think review articles fill this gap. They help, but a review captures one author's interpretation at one moment in time. Nobody updates it when new results arrive. And reviews narrate rather than decompose. A review might say "the autonomous dynamics view has been challenged," but it won't tell you that the challenge comes from two logically independent sources, neither of which has anything to do with the separate problem of preprocessing sensitivity.

Tools like Semantic Scholar and Scite.ai are better. They classify citations. They can tell you Paper A cites Paper B in a contrasting way. But they still operate at the level of documents. They cannot tell you that Paper A's claim depends on a 20 ms Gaussian smoothing kernel, that Paper C showed this choice qualitatively changes the result, and that Paper D's disagreement is about theoretical interpretation rather than data.

The citation graph itself makes things worse. It treats a supportive citation, a methodological critique, and a background nod as the same type of edge. Serra-Garcia and Gneezy [11] found that nonreplicable papers accumulate far more citations than replicable ones. Greenberg [12] traced a single biomedical claim through 242 papers and found that 94% of citations to primary data went to supportive studies. The citation graph is not neutral infrastructure. It actively amplifies unreliable claims. This is especially dangerous when median statistical power in neuroscience is around 21% [13].

**What is missing** - Every major literature platform indexes documents. None index claim structure. The missing piece is a structured claim layer tied to an explicit dependency layer. The unit would not be the paper. It would be the claim, with typed edges for support, dependency, challenge, replication, and scope restriction. A dependency layer would record what each claim actually rests on. Data source, preprocessing pipeline, model class, statistical test, species, task design. Later work could then challenge a specific dependency, replicate under altered conditions, or leave the core claim untouched while narrowing its domain of validity.

Notice what this makes visible. A study that narrows the scope of a famous claim would create a concrete restriction in the record, rather than vanishing into prose that nobody reads. A null result would become a real object, an update about where a claim fails to hold. Convergence across methods, species, or task designs would be easier to detect because the evidence would no longer be trapped inside separate narratives that no single person has time to read.

**What it would actually do** - To make this concrete, consider the rotational dynamics case. Churchland et al. 2012 would not appear as a single entry. It would be decomposed into constituent claims, each carrying its dependencies. One claim is that reaching activity shows rotational structure under a particular analysis pipeline. The dependencies include Utah-array recordings in macaques, condition averaging, Gaussian smoothing, dimensionality reduction via PCA, and jPCA. Elsayed and Cunningham's paper would then attach not to the Churchland document in general but to one specific dependency node. The detection of rotational structure depends strongly on preprocessing. That update would propagate selectively. Studies whose inferences depend on similar smoothing would be flagged. Studies reaching related conclusions through optogenetic perturbation or feedback-control modeling would remain unaffected.

This selective propagation is the core of the idea. Right now the field has two modes. Either nothing happens when a critique is published, or some individual reader happens to hold the entire web of dependencies in their head. A claim-dependency graph creates a third option, one where challenges move along actual epistemic connections rather than relying on social memory.

The implications go beyond catching errors. If five labs have tested a claim using different preprocessing, species, and task designs, that convergence is far stronger evidence than any single paper, but no existing system makes it visible. Follow-up work that tightens the scope of a claim has almost no career value today, because it is not "new." In a claim graph, it strengthens an edge, and that contribution becomes legible. Null results gain a function. A study that tests a claim under new conditions and finds it doesn't hold is not a failure but a scope restriction, represented as a real object in the graph. And claims that cannot both be true, because they assume incompatible preprocessing or theoretical commitments, would be flagged by the structure of the graph itself.

This is what I mean by a "lean for science." Not a proof assistant. Empirical science is far too messy for formal verification. But a system where claims must be situated relative to other claims, and where contradictions become visible by default rather than buried by convention.

Building this requires solving a genuinely hard design problem. What counts as a "claim"? Churchland et al. 2012 contains at least three separable claims, and reasonable people will decompose it differently. Too coarse and you lose the dependency structure. Too fine and experts can't verify entries in reasonable time. The right granularity has to be found empirically.

The architecture I'd propose has three layers. Language models extract candidate claims, dependencies, and edges from paper text. A graph-level algorithm propagates sensitivity flags along dependency chains, detects inconsistencies between claims with incompatible assumptions, and surfaces clusters of claims that converge through independent paths. Domain experts adjudicate the outputs of both layers. The schema would cover a fixed set of dependency types (data source, species, task, preprocessing, statistical test, model class, theoretical commitment) and edge types (supports, depends-on, challenges, restricts-scope, replicates). Whether this is expressive enough is one of the things the pilot would test.

**The experiment** - Three things have recently changed that make this buildable now.

First, standardized neural data archives mean claims can link back to queryable data rather than figure panels. DANDI now hosts over 1,000 datasets. The Neural Latents Benchmark provides preprocessed population recordings from the exact motor cortex experiments at issue.

Second, language models have dropped the cost of structured extraction from scientific text by roughly an order of magnitude. What previously required expert annotators reading each paper can now be bootstrapped by LLM extraction with expert correction. That makes a 300–500 paper corpus tractable for a small team.

Third, AI-assisted writing is accelerating paper production without adding any structure [14]. The problem is getting worse faster than the field is building tools to address it.

I propose building a claim-dependency graph for a single, well-documented subfield. Motor neuroscience is the case I've developed here, but the approach is not specific to it. Assemble a corpus of 300 to 500 core papers. Define the schema. Use language models to propose candidate entries. Have domain experts adjudicate. Release the graph publicly.

The most informative test is a backtest. Build the graph using only papers published before some cutoff, say 2018. Then ask whether the dependency structure flags problems the field took years to notice on its own. Does it surface the tension between autonomous and input-driven dynamics before Sauerbrei 2020? The possibility that reaching results might not generalize to grasping before Suresh 2020? The fact that Elsayed and Cunningham's preprocessing concerns should have prompted more caution about downstream claims? If yes, the infrastructure would have saved years of confused debate. If no, either the schema is too coarse or the dependencies that matter were not visible in the earlier literature, and the bottleneck is elsewhere.

No existing institution is built to maintain this. Universities don't reward infrastructure work. Grant panels evaluate novelty. Publishers have no reason to flag contradictions in their own product. It would require a dedicated team, probably something like a Focused Research Organization.

I want to be honest about how this could fail. The schema may be too rigid. The graph may capture real structure but fail the backtest, surfacing only tensions already obvious from the pre-2018 literature. Or the backtest may succeed for known cases but fail to generalize, because we built it knowing what to look for. Any of these outcomes would be worth knowing.

Motor neuroscience is a good pilot domain because it has a bounded literature, standardized data formats, and active disputes with identifiable dependency structure. But it is not unique. Gene Ontology gave biology a shared vocabulary. The Protein Data Bank gave structural biology a shared archive. Lean gave mathematics a shared verification layer. Neuroscience has already built the data archives and the computation standards. What it still lacks is the glue, a structured and queryable layer connecting claims to evidence to data to code. Building that layer for one subfield is a tractable experiment. If it works, it becomes a template.

## References

1. A. P. Georgopoulos, A. B. Schwartz, R. E. Kettner, "Neuronal population coding of movement direction," *Science*, vol. 233, pp. 1416–1419, 1986.
2. E. Todorov, "Direct cortical control of muscle activation in voluntary arm movements," *Nature Neuroscience*, vol. 3, pp. 391–398, 2000.
3. T. D. Sanger, "Theoretical considerations for the analysis of population coding in motor cortex," *Neural Computation*, vol. 6, pp. 29–37, 1994.
4. A. Ramot, F. H. Taschbach, Y. C. Yang, et al., "Motor learning refines thalamic influence on motor cortex," *Nature*, 2025.
5. S. H. Scott, "Optimal feedback control and the neural basis of volitional motor control," *Nature Reviews Neuroscience*, vol. 5, pp. 532–546, 2004.
6. M. M. Churchland, J. P. Cunningham, et al., "Neural population dynamics during reaching," *Nature*, vol. 487, pp. 51–56, 2012.
7. B. A. Sauerbrei, J.-Z. Guo, et al., "Cortical pattern generation during dexterous movement is input-driven," *Nature*, vol. 577, pp. 386–391, 2020.
8. H. T. Kalidindi et al., "Rotational dynamics in motor cortex are consistent with a feedback controller," *eLife*, vol. 10, e67256, 2021.
9. G. F. Elsayed, J. P. Cunningham, "Structure in neural population recordings: an expected byproduct of simpler phenomena?" *Nature Neuroscience*, vol. 20, pp. 1310–1318, 2017.
10. A. K. Suresh, J. M. Goodman, et al., "Neural population dynamics in motor cortex are different for reach and grasp," *eLife*, vol. 9, e58848, 2020.
11. M. Serra-Garcia, U. Gneezy, "Nonreplicable publications are cited more than replicable ones," *Science Advances*, vol. 7, eabd1705, 2021.
12. S. A. Greenberg, "How citation distortions create unfounded authority: analysis of a citation network," *BMJ*, vol. 339, b2680, 2009.
13. K. S. Button et al., "Power failure: why small sample size undermines the reliability of neuroscience," *Nature Reviews Neuroscience*, vol. 14, pp. 365–376, 2013.
14. W. Liang, Z. Izzo, Y. Zhang, et al., "Monitoring AI-modified content at scale: a case study on the impact of ChatGPT on AI conference peer reviews," *Proceedings of the 41st International Conference on Machine Learning*, PMLR 235, pp. 29575–29620, 2024.