Large Language Models as Cognitive Models

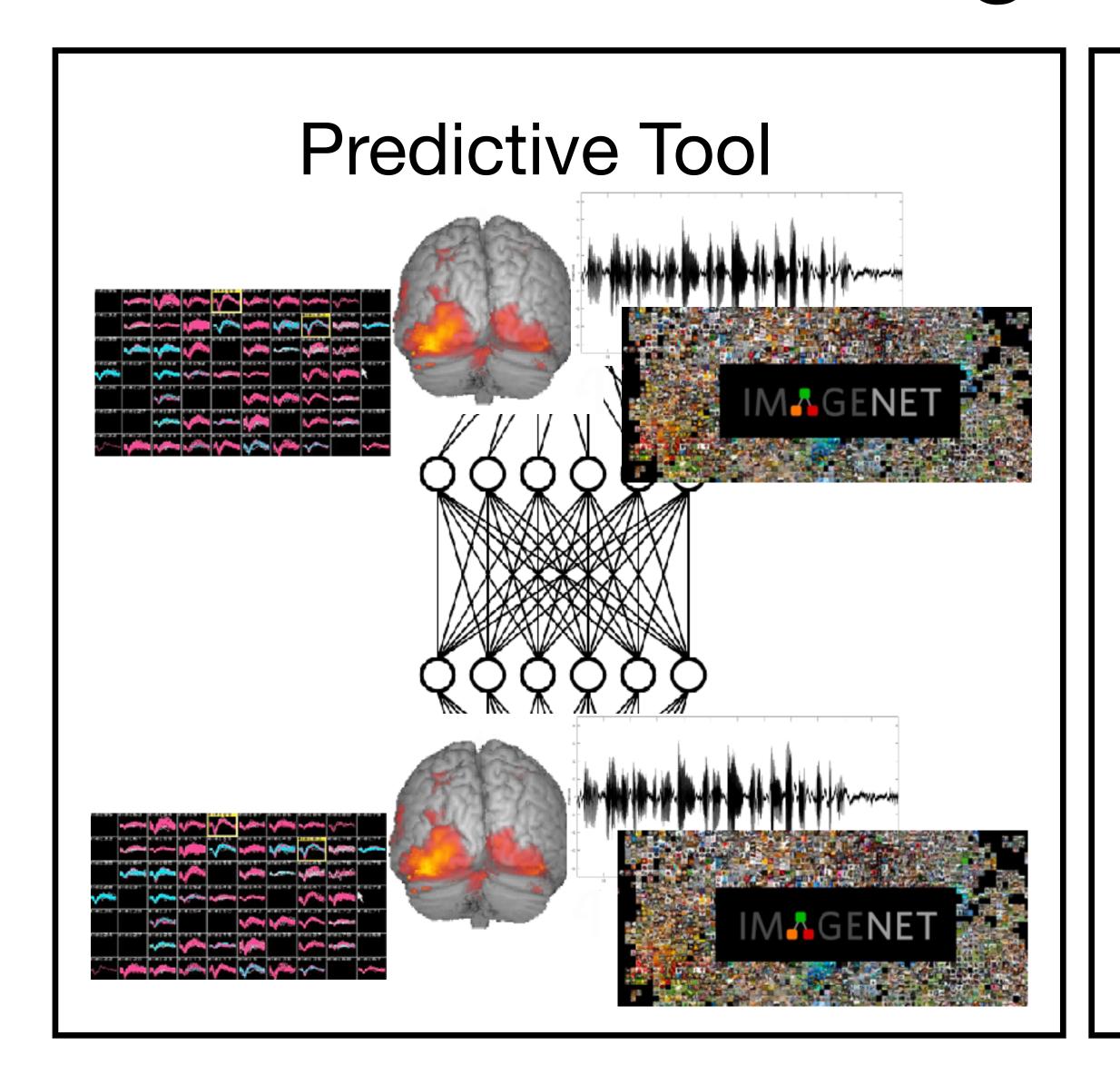
Neuroscience Forum Workshop: Exploring the Bidirectional Relationship Between Artificial Intelligence and Neuroscience



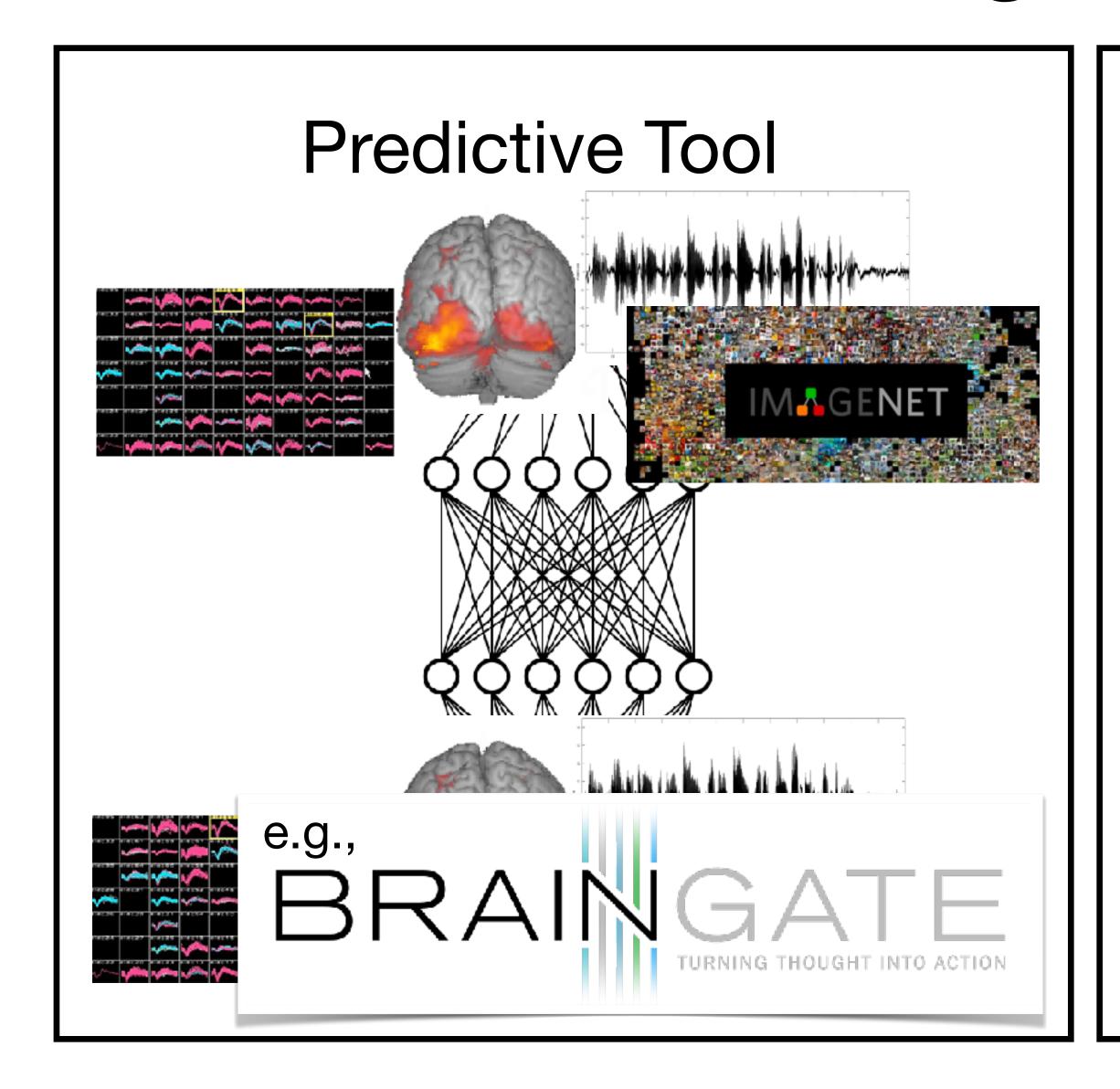
Disclosures

- Current/Recent Research Support: NSF, DARPA, IARPA, Google, ONR, NIH
- Consulting/Speaking Engagements: Google Deepmind (Current Employee);
 Signify (One-Time Consultant)

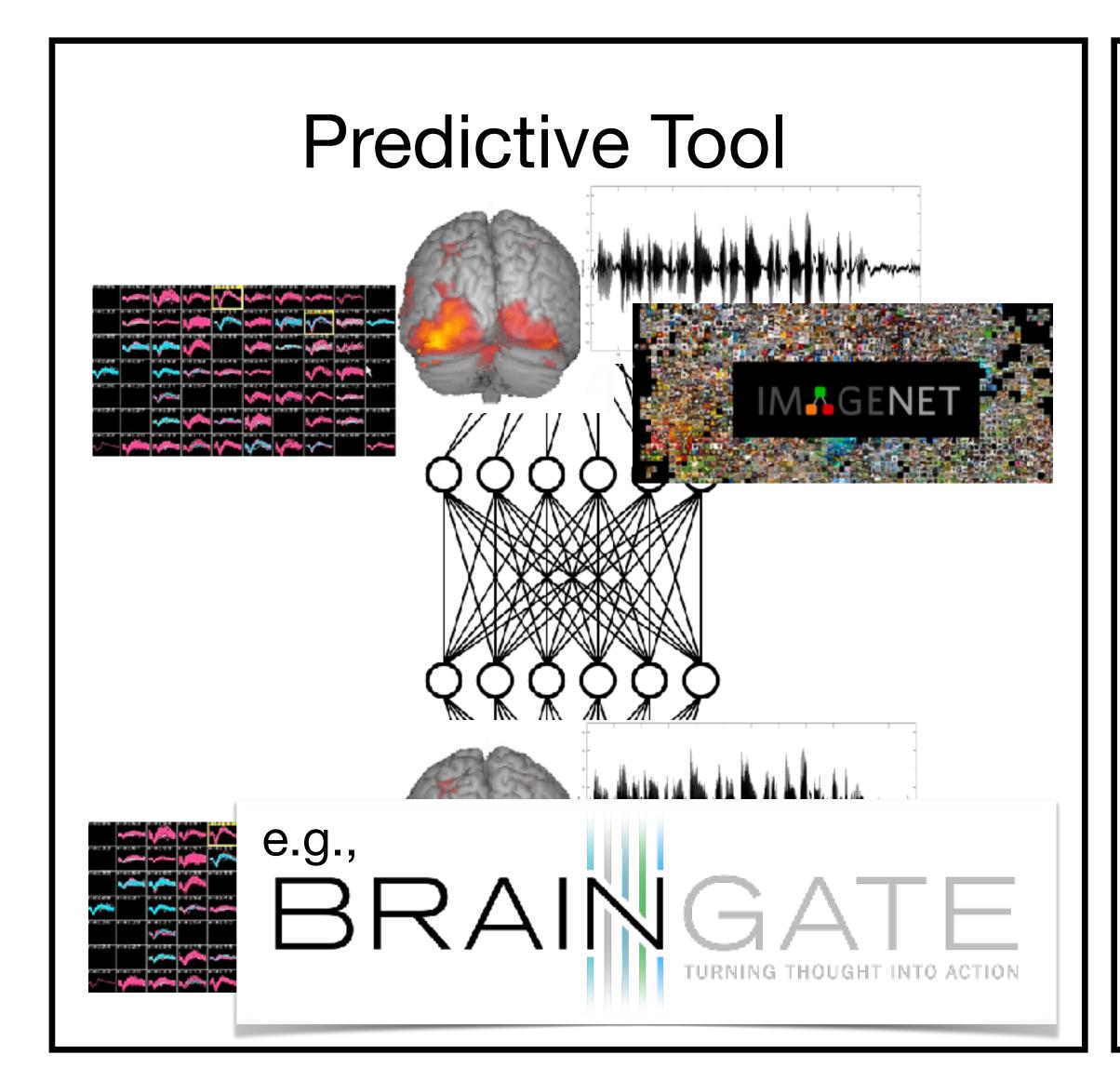
Predictive Tool **Explanatory Model**

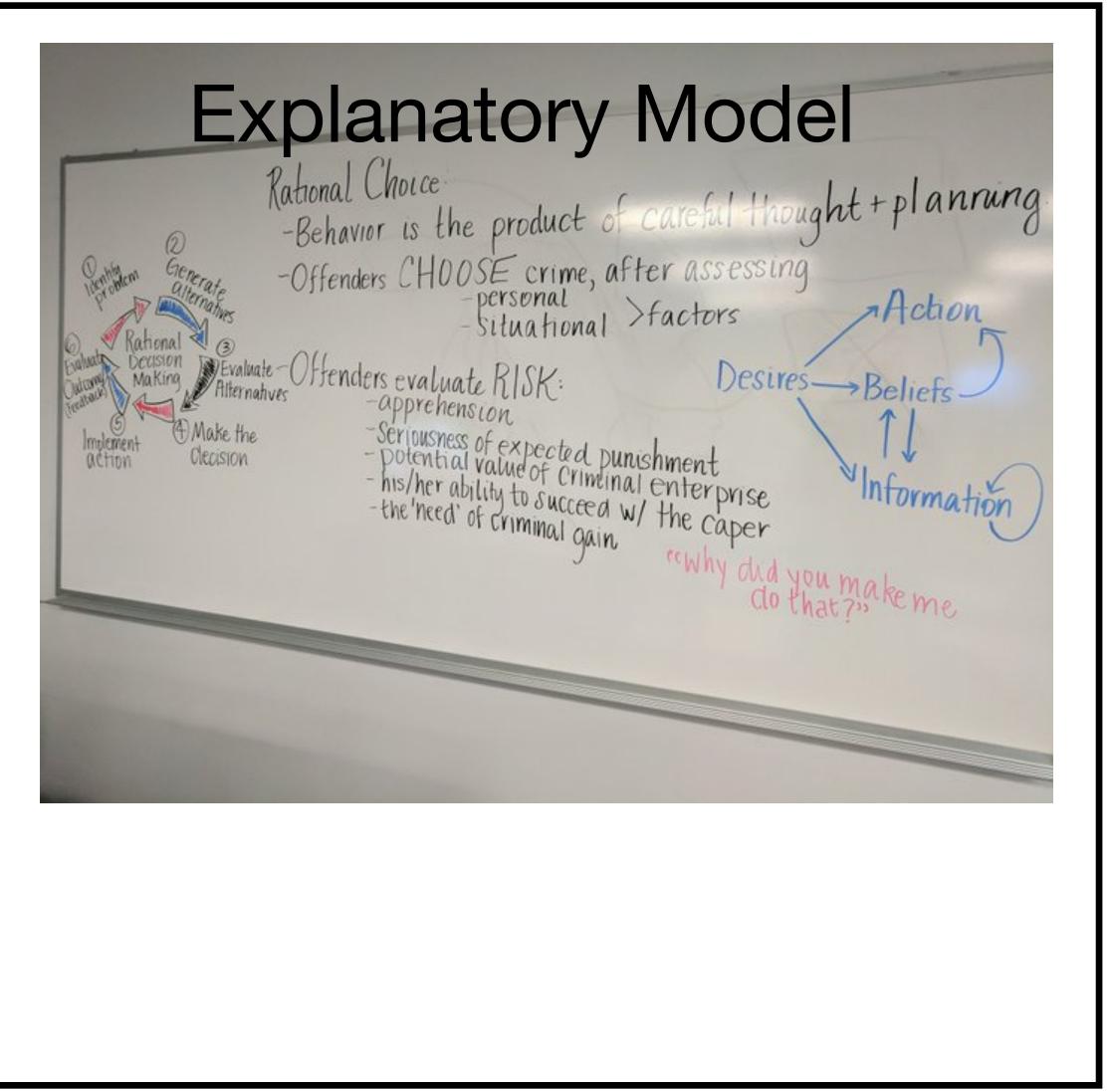


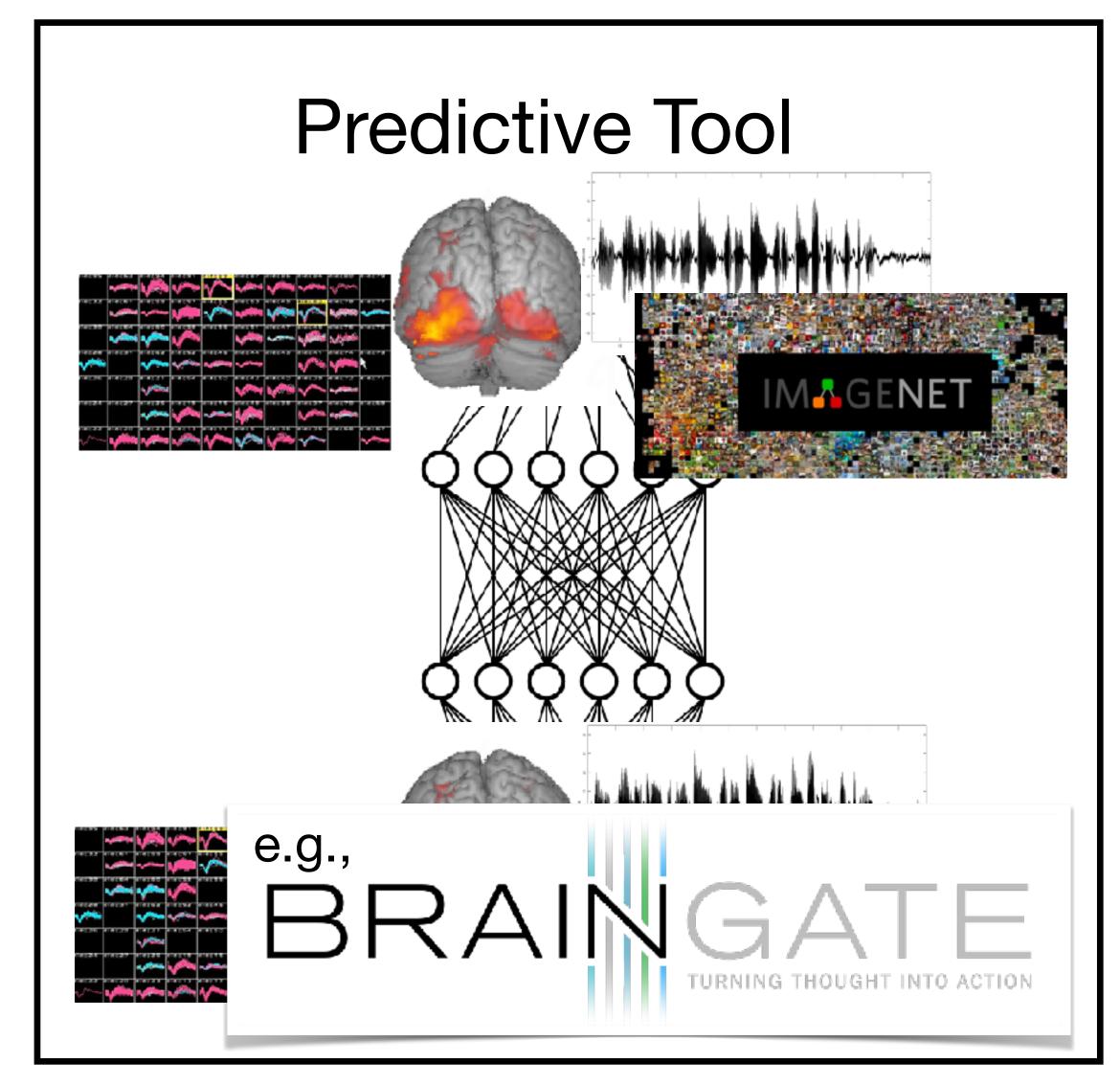
Explanatory Model

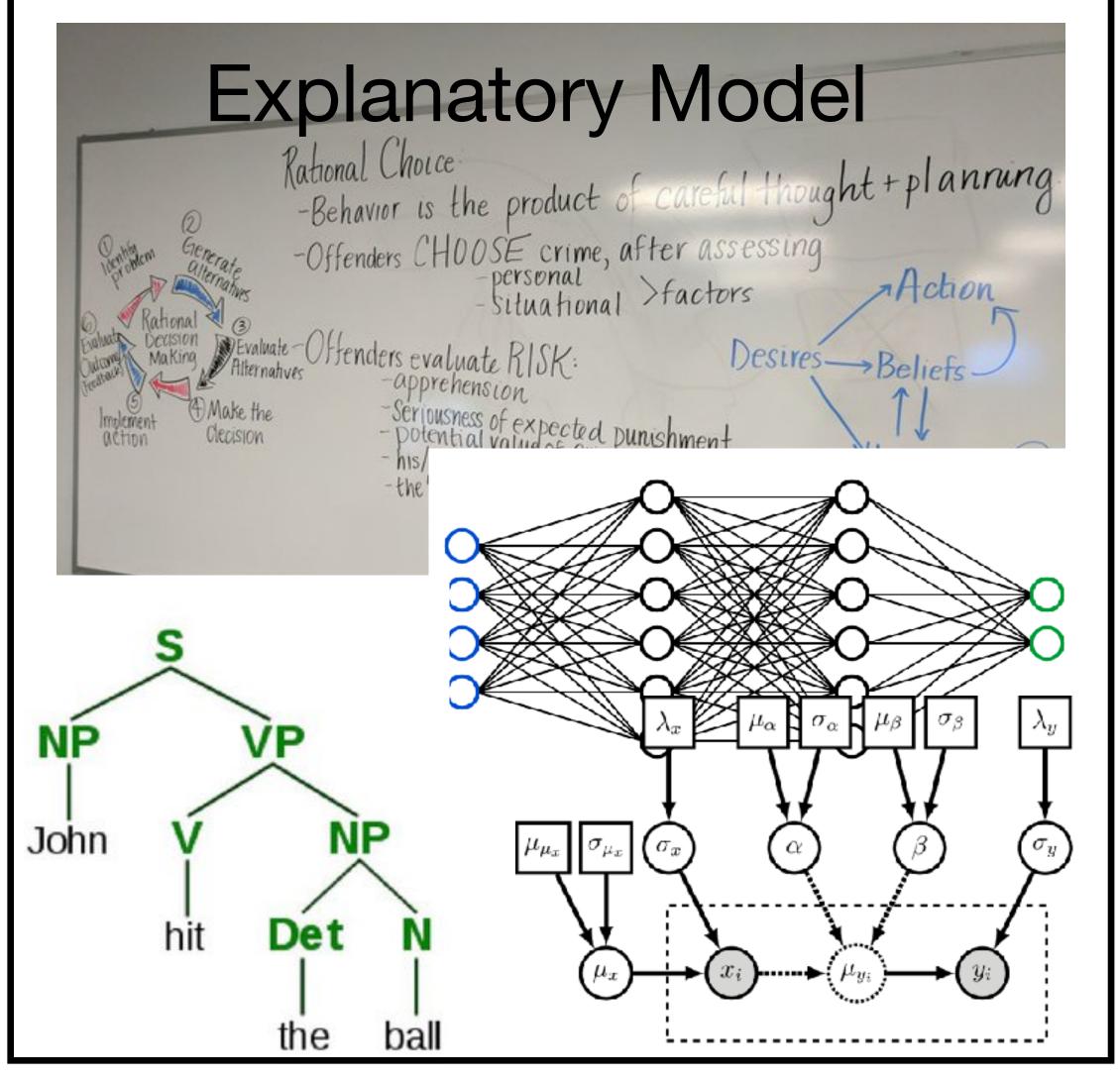


Explanatory Model

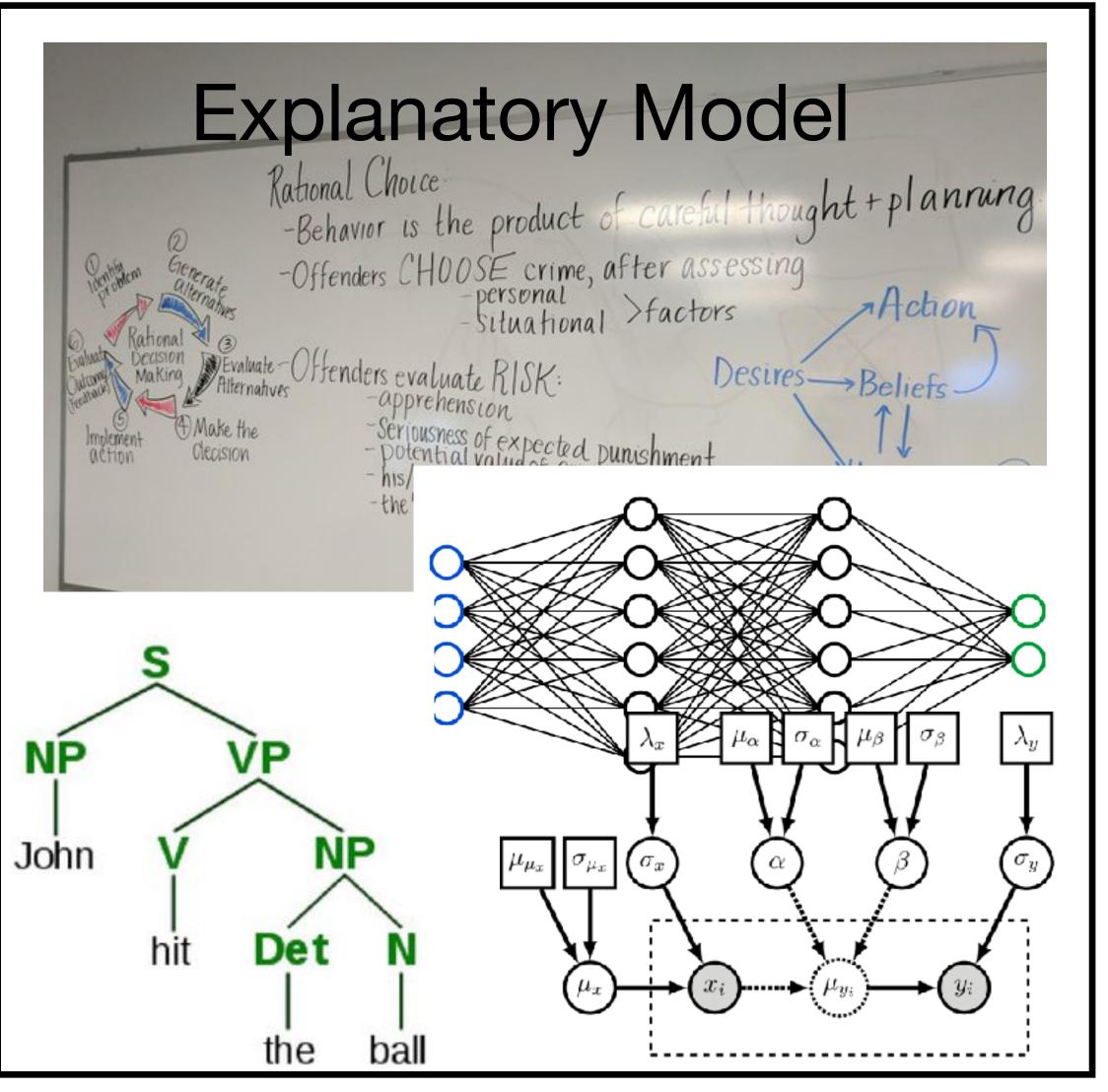












Reasons not to use LLMs as Cognitive Models

Reasons *not* to use LLMs as Cognitive Models

LLMs were not designed as cognitive models.

The success was "stumbled upon" via good engineering, but they do not represent any coherent theory. And, they still are woefully un-human-like in so many ways.

Reasons not to use LLMs as Cognitive Models

LLMs were not designed as cognitive models.

The success was "stumbled upon" via good engineering, but they do not represent any coherent theory. And, they still are woefully un-human-like in so many ways.

LLMs are black boxes.

We do not understand how they work. We cannot simply replace one black box—i.e., the human brain and mind—with another black box—i.e, the artificial neural network.

Reasons to use LLMs as Cognitive Models

LLMs were not designed as cognitive models.

Yes. But their ability to simulate human behavior across so many varied domains is worth taking seriously. And in certain some cases (namely, language) they represent our best computational model of human behavior by a long shot.

LLMs are black boxes.

Reasons to use LLMs as Cognitive Models

LLMs were not designed as cognitive models.

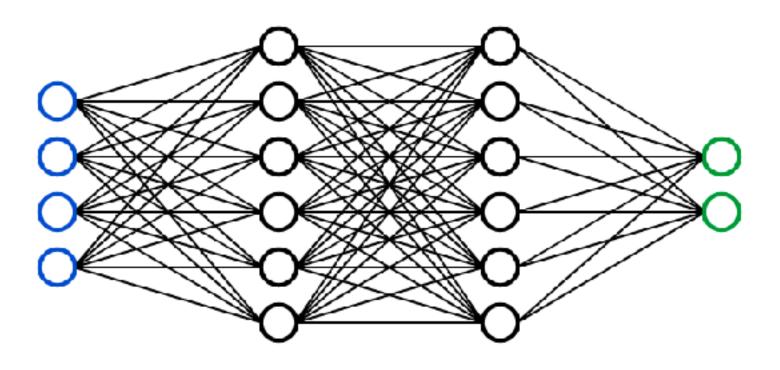
Yes. But their ability to simulate human behavior across so many varied domains is worth taking seriously. And in certain some cases (namely, language) they represent our best computational model of human behavior by a long shot.

LLMs are black boxes.

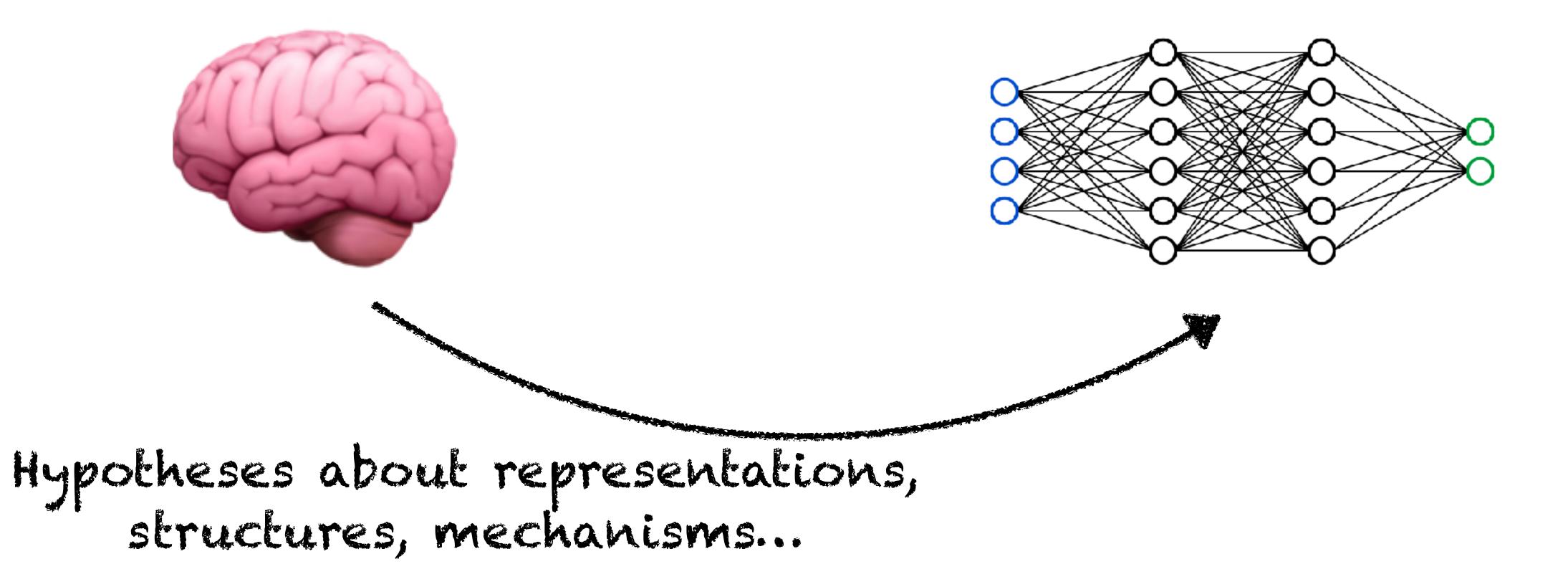
Not necessarily. In fact, recent work shows LLMs might be more interpretable than we often assume, and uncovering these mechanisms could form a virtuous cycle that advances cognitive neuroscience+AI in tandem.

Virtuous Cycle between (Generative) Al and Cognitive Neuroscience

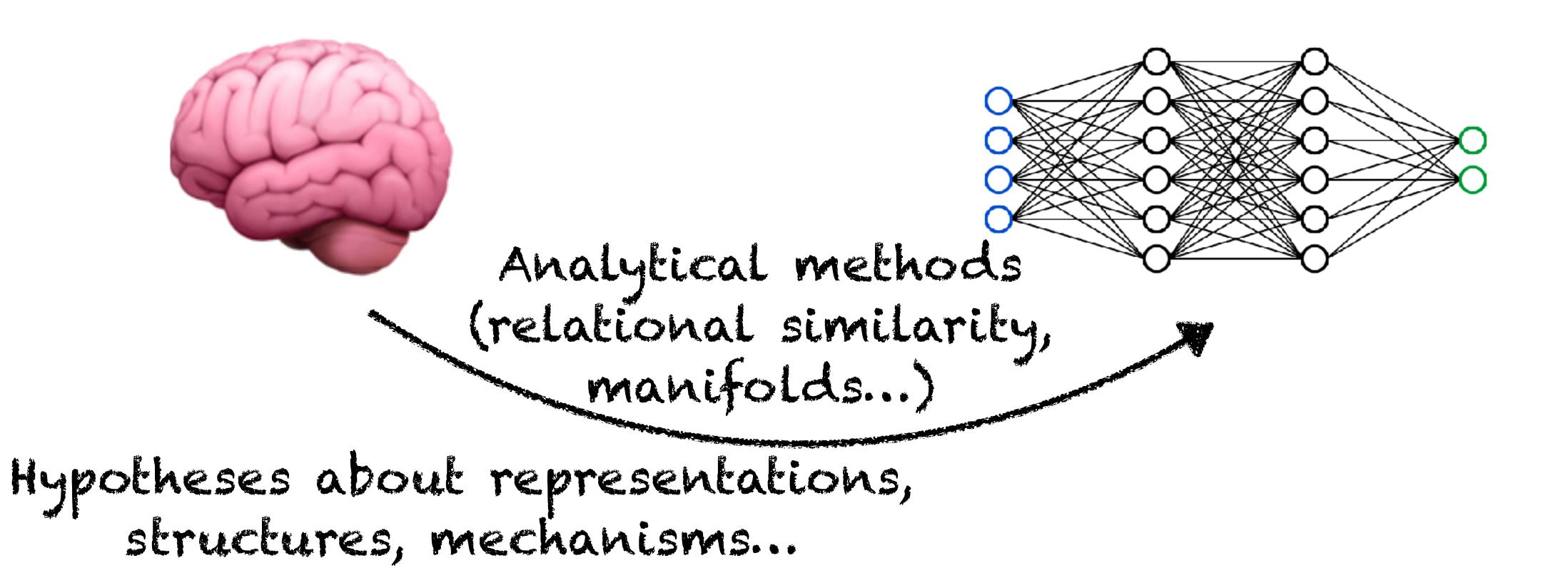




Virtuous Cycle between (Generative) Al and Cognitive Neuroscience



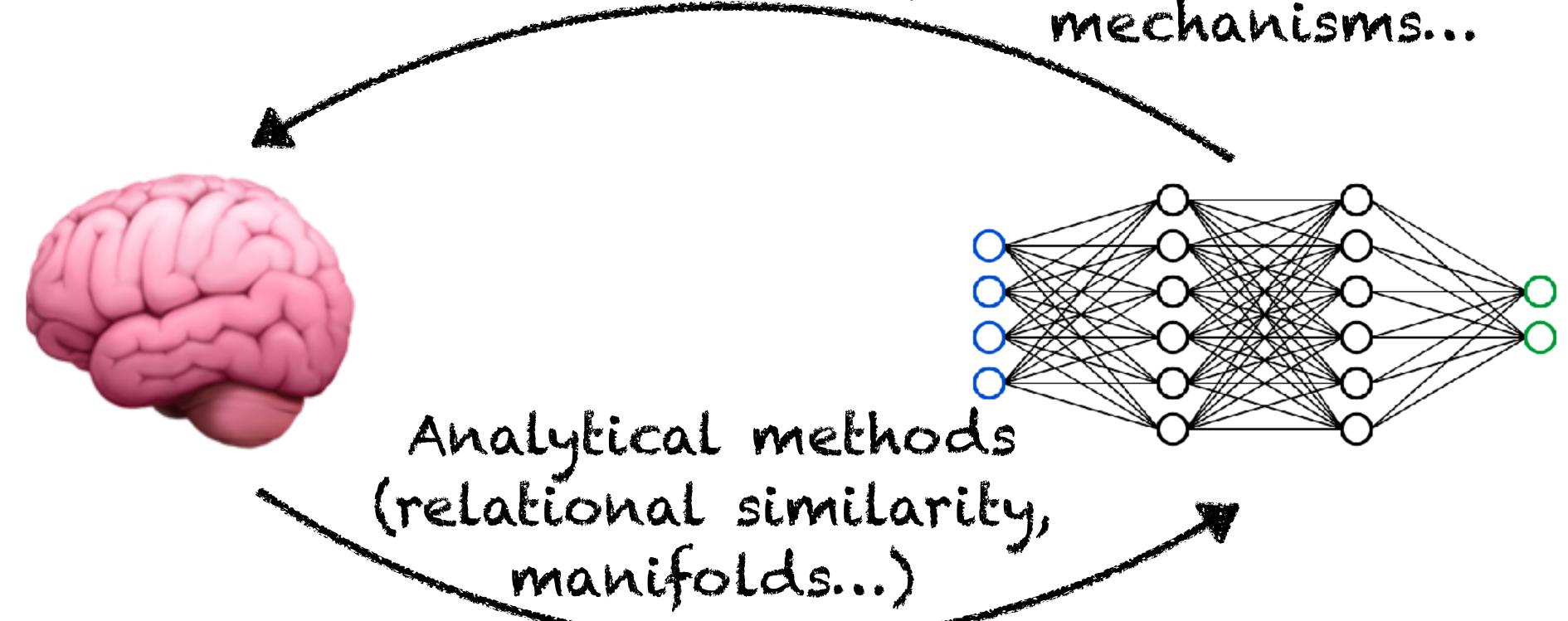
Virtuous Cycle between (Generative) Al and Cognitive Neuroscience



Virtuous Cycle between (Generative) Al and Cognitive Neuroscience

Hypotheses about

Hypotheses about representations, structures, mechanisms...

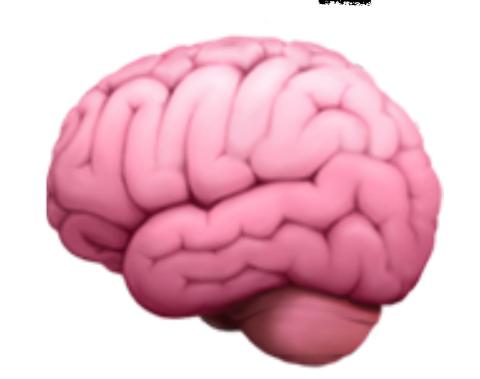


Hypotheses about representations, structures, mechanisms...

Virtuous Cycle between (Generative) Al and Cognitive Neuroscience

Hypotheses about

Forward-engineering, causal interventions ("lesions")



Hypotheses about representations, structures, mechanisms...

Analytical methods (relational similarity, manifolds...)

Hypotheses about representations, structures, mechanisms...

Virtuous Cycle between (Generative) Al and Cognitive Neuroscience Hypotheses about

Forward-engineering, causal interventions
("lesions")



Analytical methods (relational similarity, manifolds...)

Hypotheses about representations, structures, mechanisms...

Hypotheses about representations, structures, mechanisms...



Proofs of Concept

- LLMs are not merely "stochastic parrots". They often contain interpretable internal structure which makes use of modular functionally-specialized components.
- These components can sometimes resemble known brain mechanisms, e.g., mechanisms for cognitive control in the prefrontal cortex
- LLMs, off the shelf, exhibit behaviors for which we don't have good competing computational cognitive models, like the ability to learn flexibly both in context and in weights. Understanding these mechanisms and aligning them to those of humans could jointly offer new theories of human cognition and improve LLMs robustness, interpretability, and efficiency

Proofs of Concept

- LLMs are not merely "stochastic parrots". They often contain interpretable internal structure which makes use of modular functionally-specialized components.
- These components can sometimes resemble known brain mechanisms,
 e.g., mechanisms for cognitive control in the prefrontal cortex
- LLMs, off the shelf, exhibit behaviors for which we don't have good competing computational cognitive models, like the ability to learn flexibly both in context and in weights. Understanding these mechanisms and aligning them to those of humans could jointly offer new theories of human cognition and improve LLMs robustness, interpretability, and efficiency

What is the capital of France?

Paris

What is the capital of Poland?

Warsaw

Merullo, Jack, Carsten Eickhoff, and Ellie Pavlick. "Language models implement simple word2vec-style vector arithmetic." *NAACL* (2024).

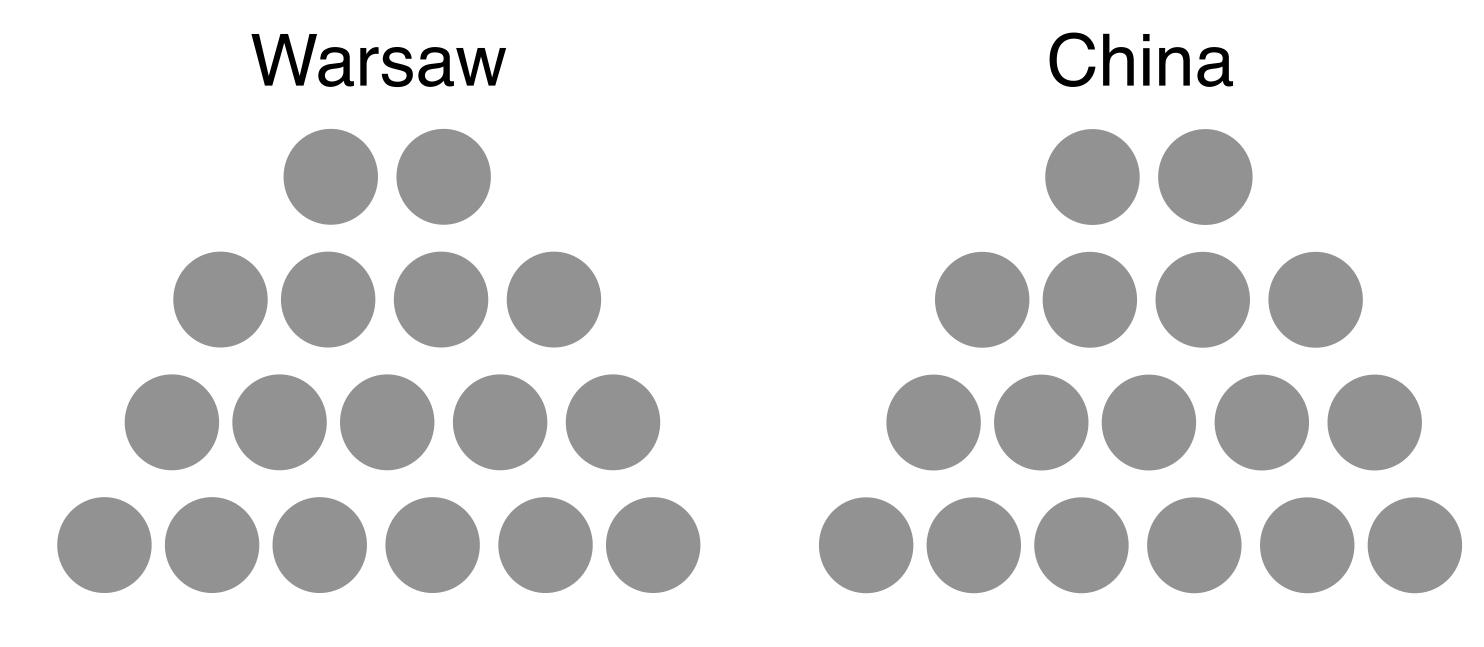
Possibility #1: Idiosyncratic

Possibility #2: Systematic

```
P(Warsaw|Poland & ...
of & capital & ...
of Poland & ...)
```

```
f|f(France) = Paris
f(Poland) = Warsaw
```

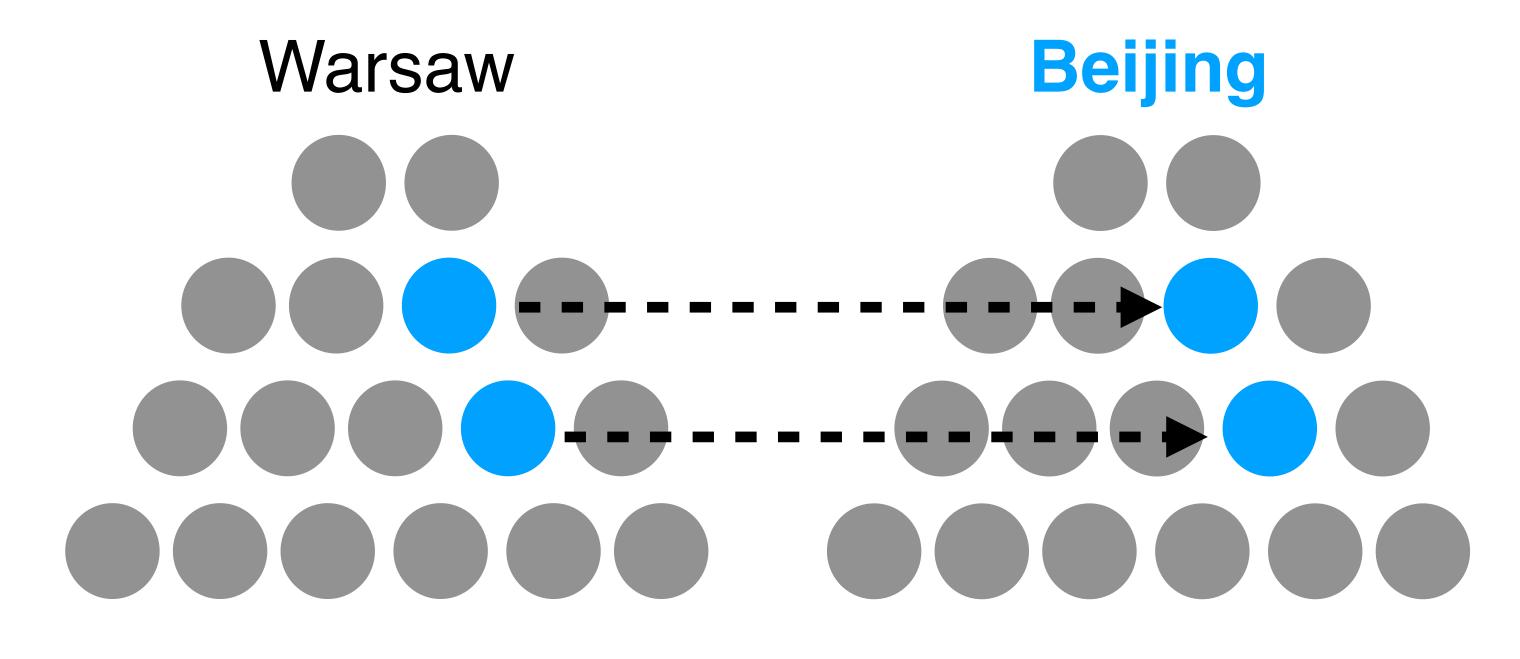
implement simple wordzvec-style vector arithmetic." NAACL (2024).

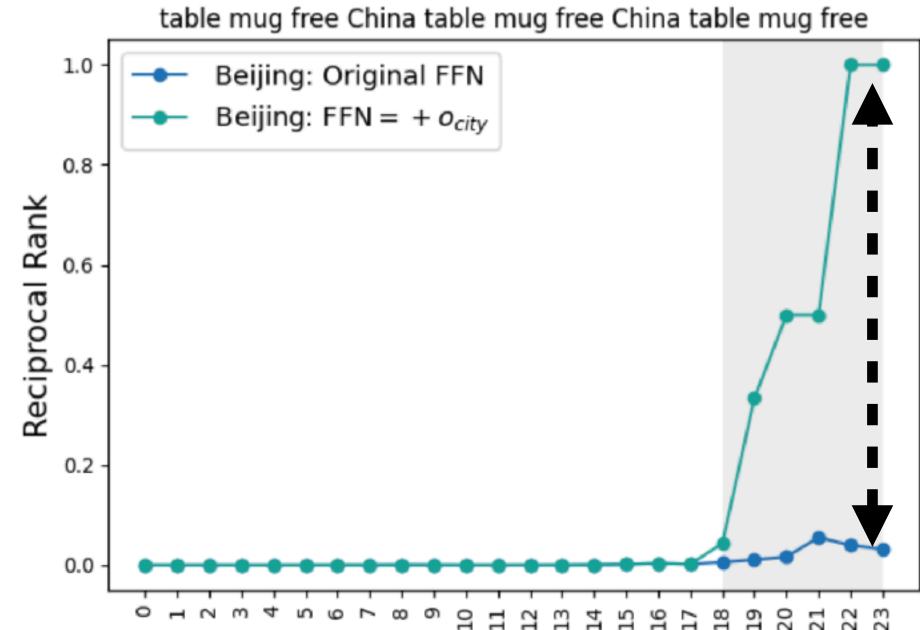


What is the capital of Poland?

beep boop China beep boop

Merullo, Jack, Carsten Eickhoff, and Ellie Pavlick. "Language models implement simple word2vec-style vector arithmetic." *NAACL* (2024).





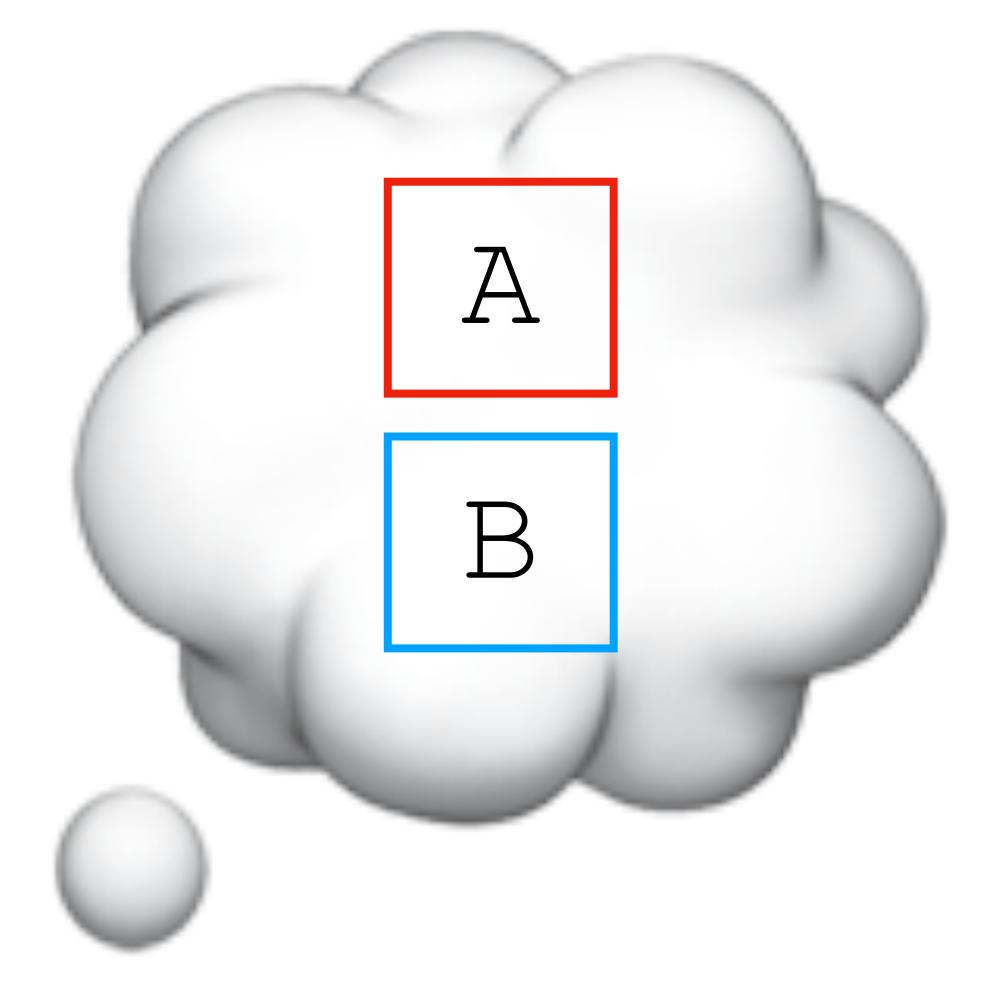
What is the capital of Poland?

beep boop China beep boop

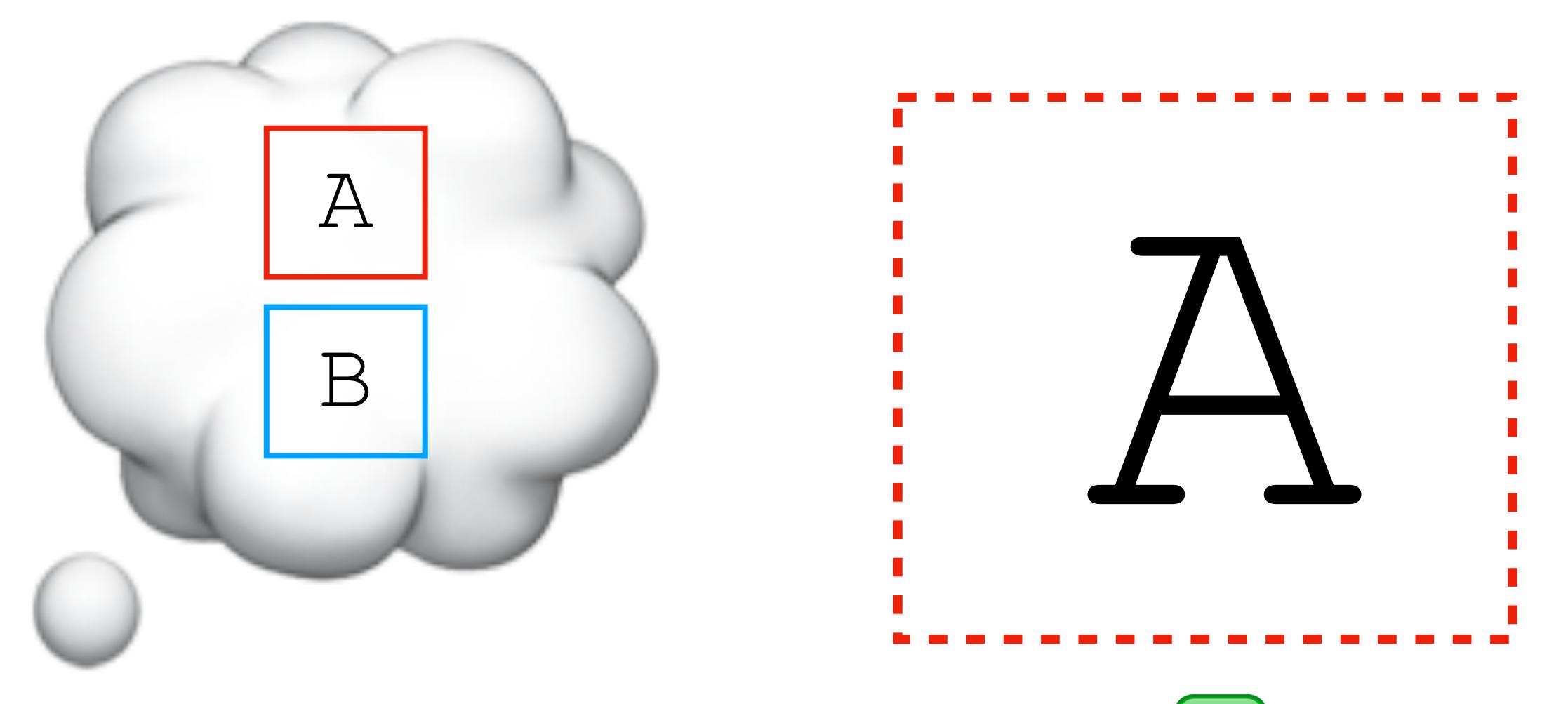
Merullo, Jack, Carsten Eickhoff, and Ellie Pavlick. "Language models implement simple word2vec-style vector arithmetic." *NAACL* (2024).

Proofs of Concept

- LLMs are not merely "stochastic parrots". They often contain interpretable internal structure which makes use of modular functionally-specialized components.
- These components can sometimes resemble known brain mechanisms, e.g., mechanisms for cognitive control in the prefrontal cortex
- LLMs, off the shelf, exhibit behaviors for which we don't have good competing computational cognitive models, like the ability to learn flexibly both in context and in weights. Understanding these mechanisms and aligning them to those of humans could jointly offer new theories of human cognition and improve LLMs robustness, interpretability, and efficiency



Traylor, Aaron, et al. "Transformer Mechanisms Mimic Frontostriatal Gating Operations When Trained on Human Working Memory Tasks." arXiv preprint arXiv:2402.08211 (2024).



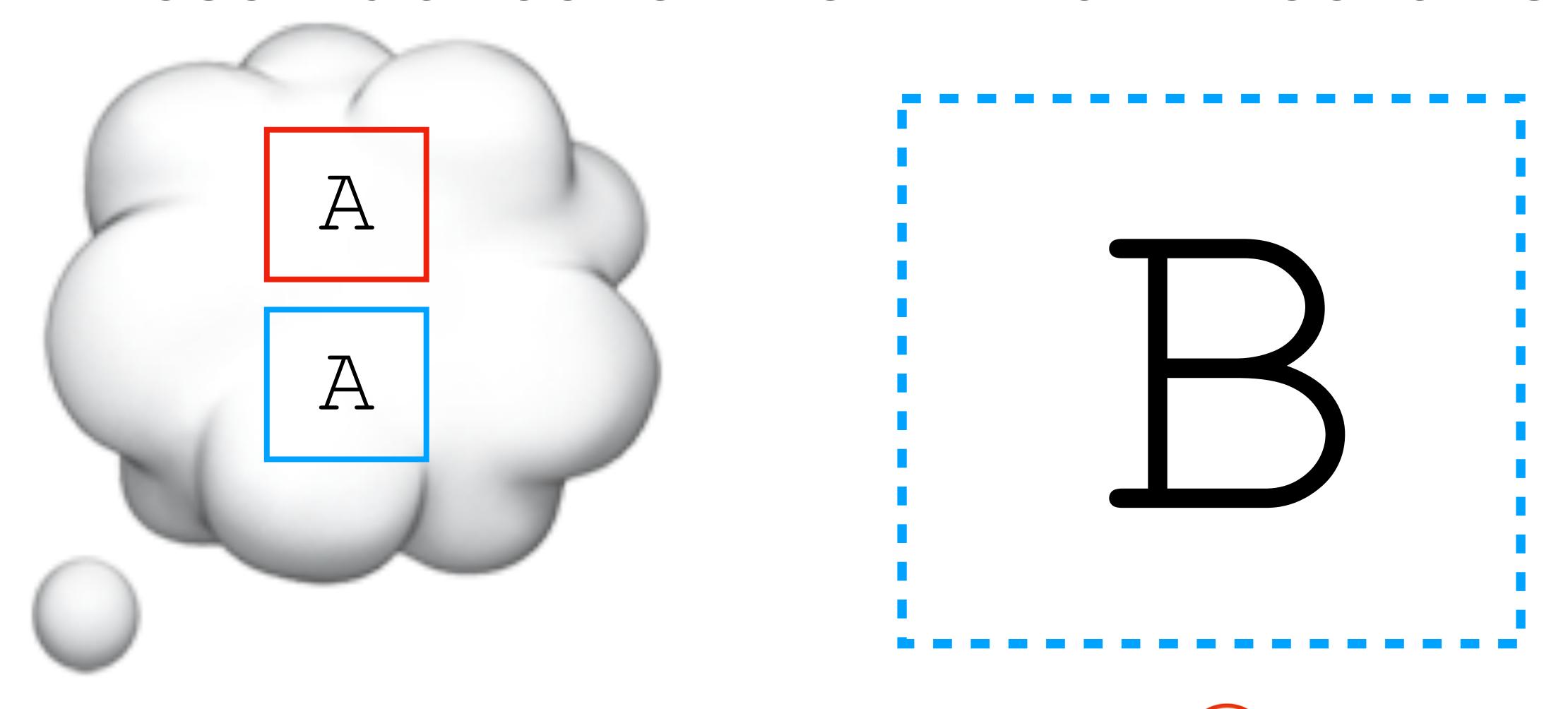
Traylor, Aaron, et al. "Transformer Mechanisms Mimic striatal Gating Operations When Trained on Human Working Memory Tasks." arXiv preprint arXiv:2402.08211 (2024).



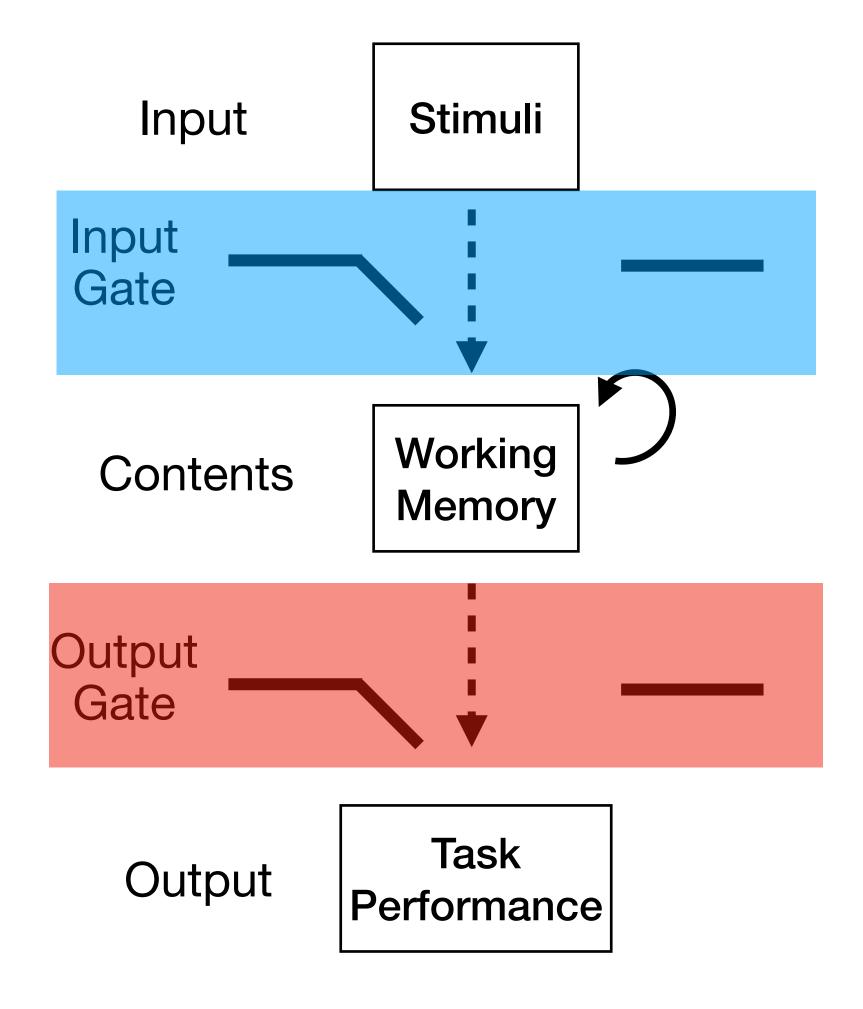
Traylor, Aaron, et al. "Transformer Mechanisms Mimic Francostriatal Gating Operations When Trained on Human Working Memory Tasks." arXiv preprint arXiv:2402.08211 (2024).



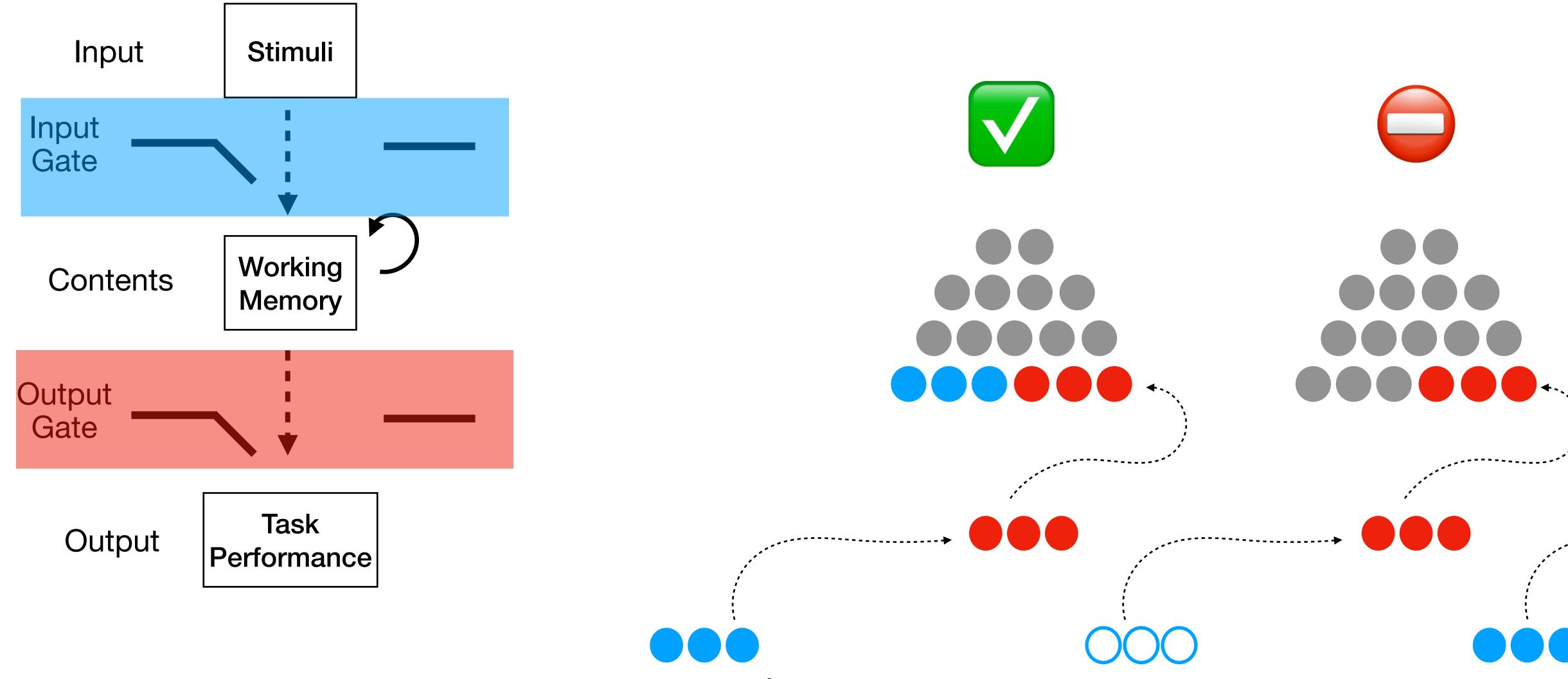
Traylor, Aaron, et al. "Transformer Mechanisms Mimic Floriostriatal Gating Operations When Trained on Human Working Memory Tasks." arXiv preprint arXiv:2402.08211 (2024).



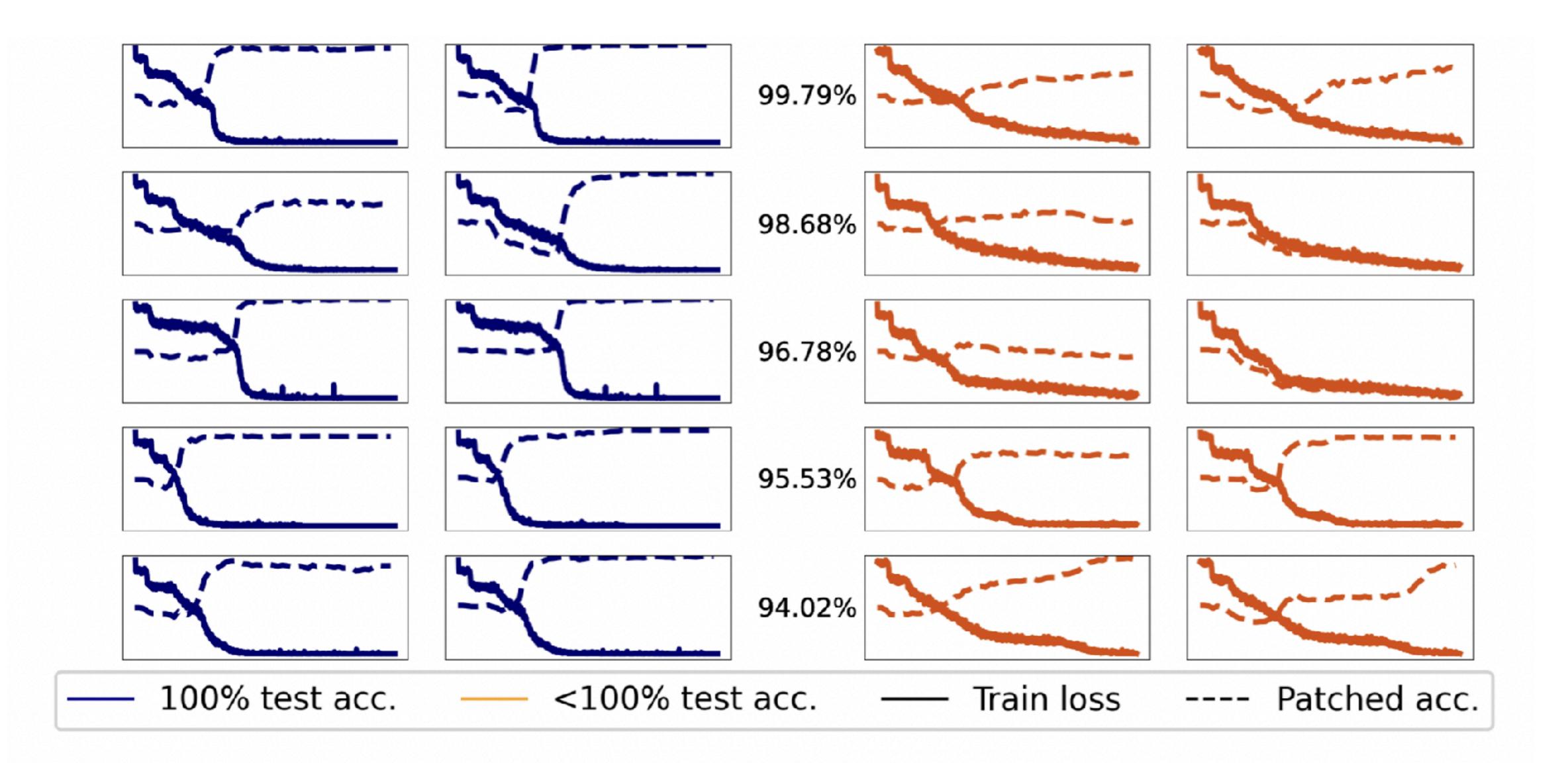
Traylor, Aaron, et al. "Transformer Mechanisms Mimic Francostriatal Gating Operations When Trained on Human Working Memory Tasks." arXiv preprint arXiv:2402.08211 (2024).

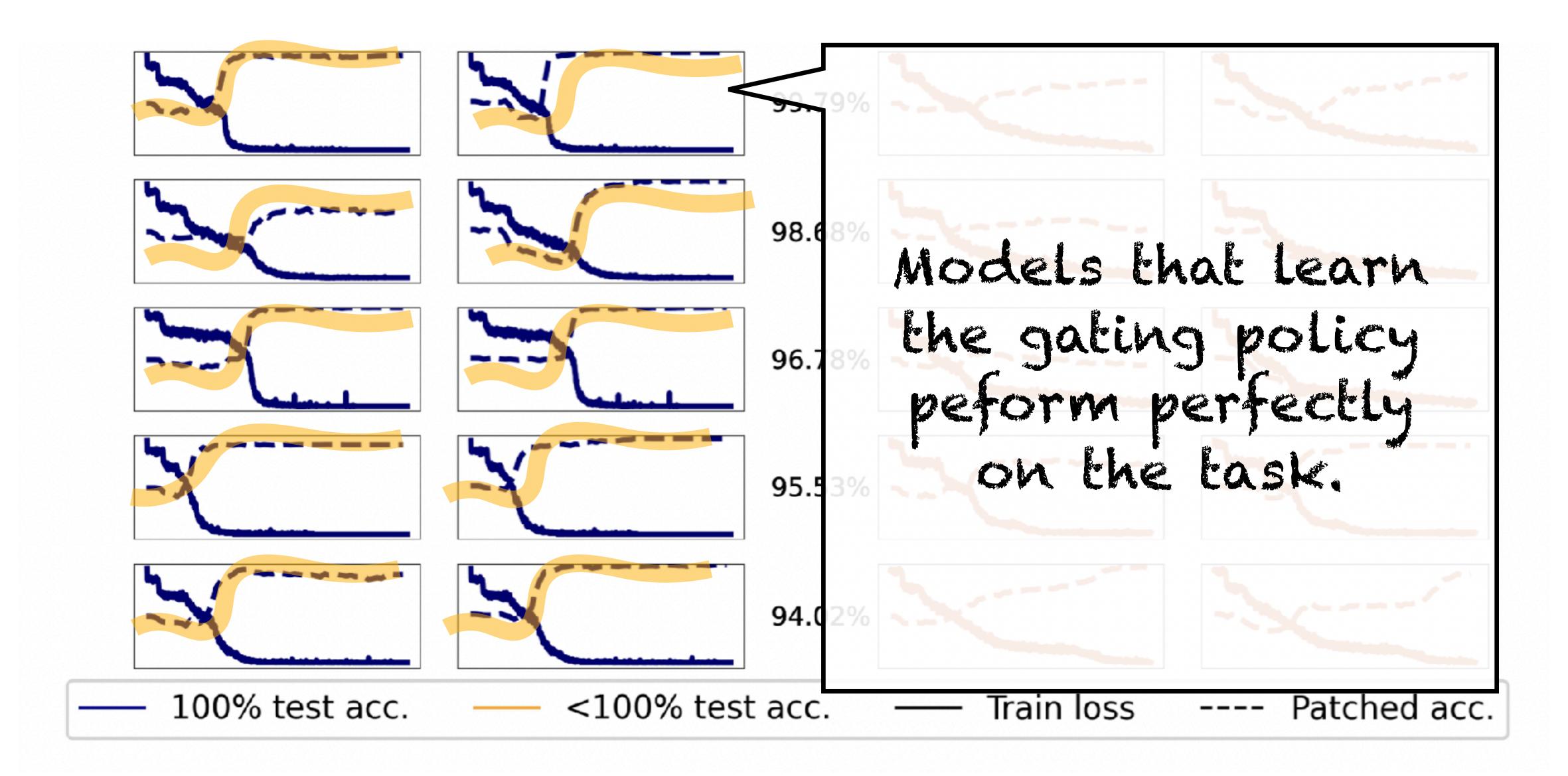


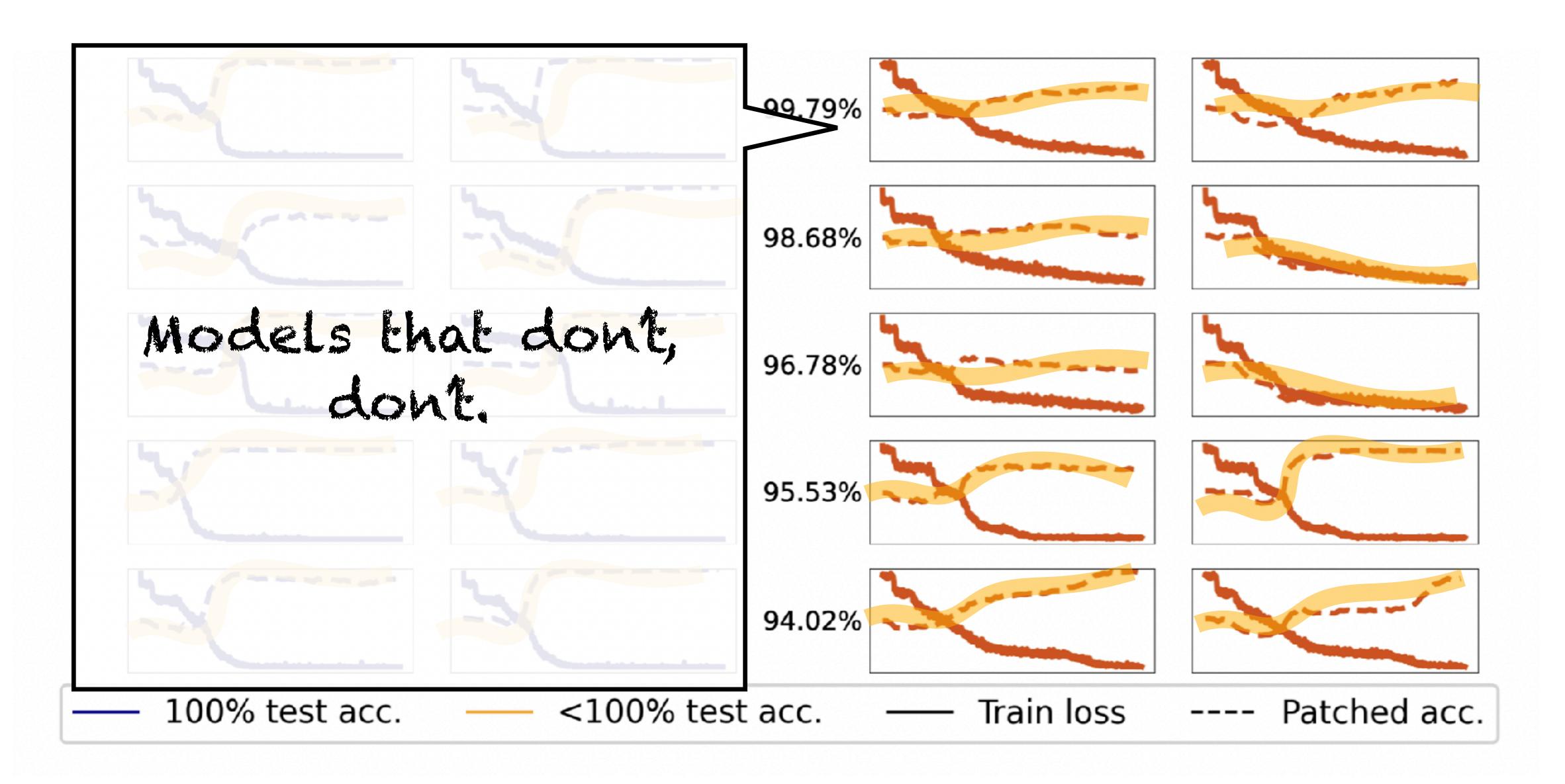
Traylor, Aaron, et al. "Transformer Mechanisms Mimic Frontostriatal Gating Operations When Trained on Human Working Memory Tasks." arXiv preprint arXiv:2402.08211 (2024).



store blue A ignore blue B store blue B

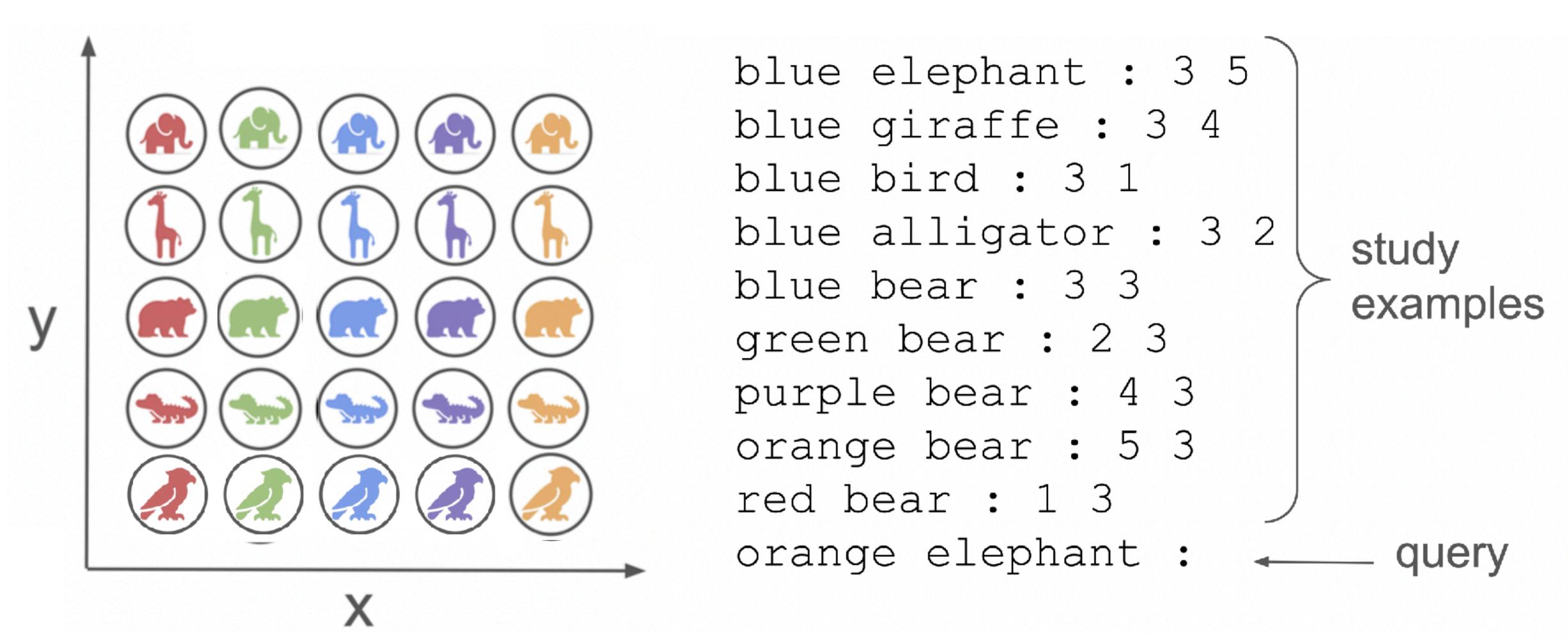




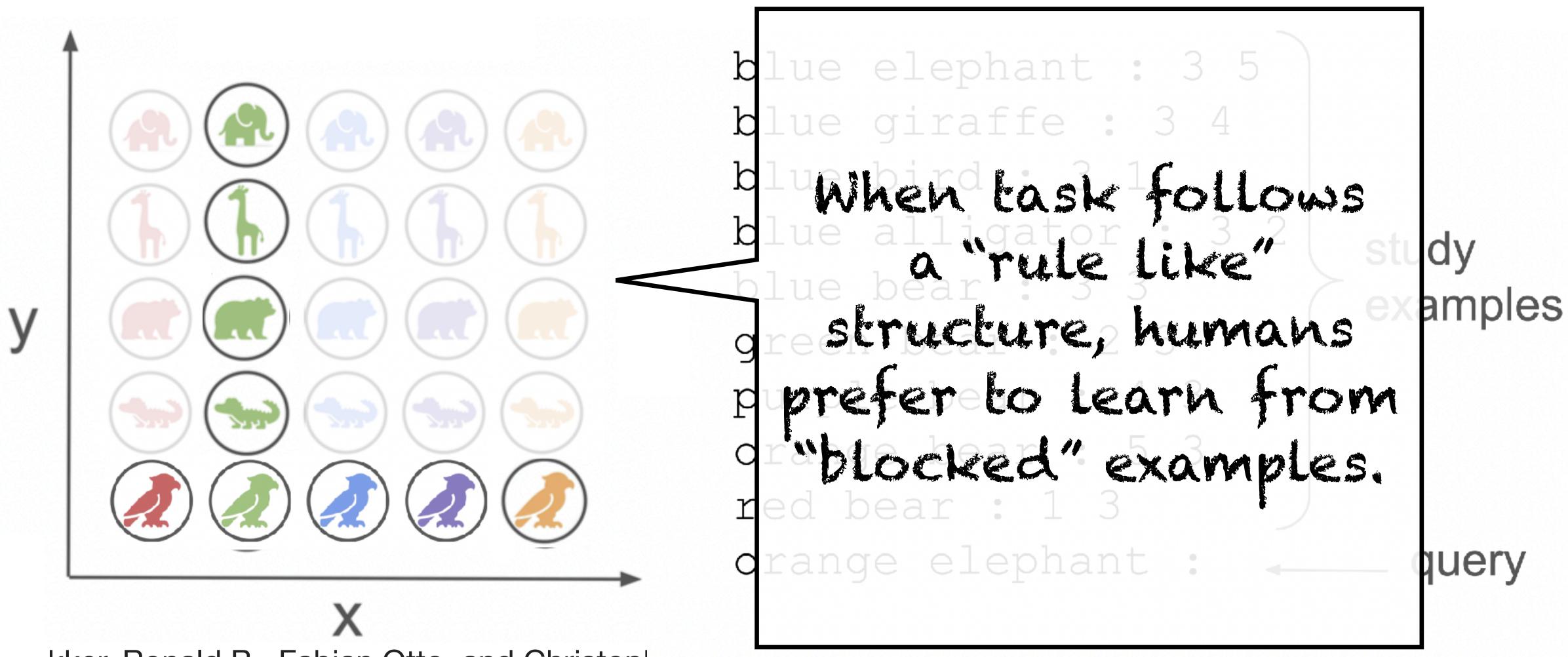


Proofs of Concept

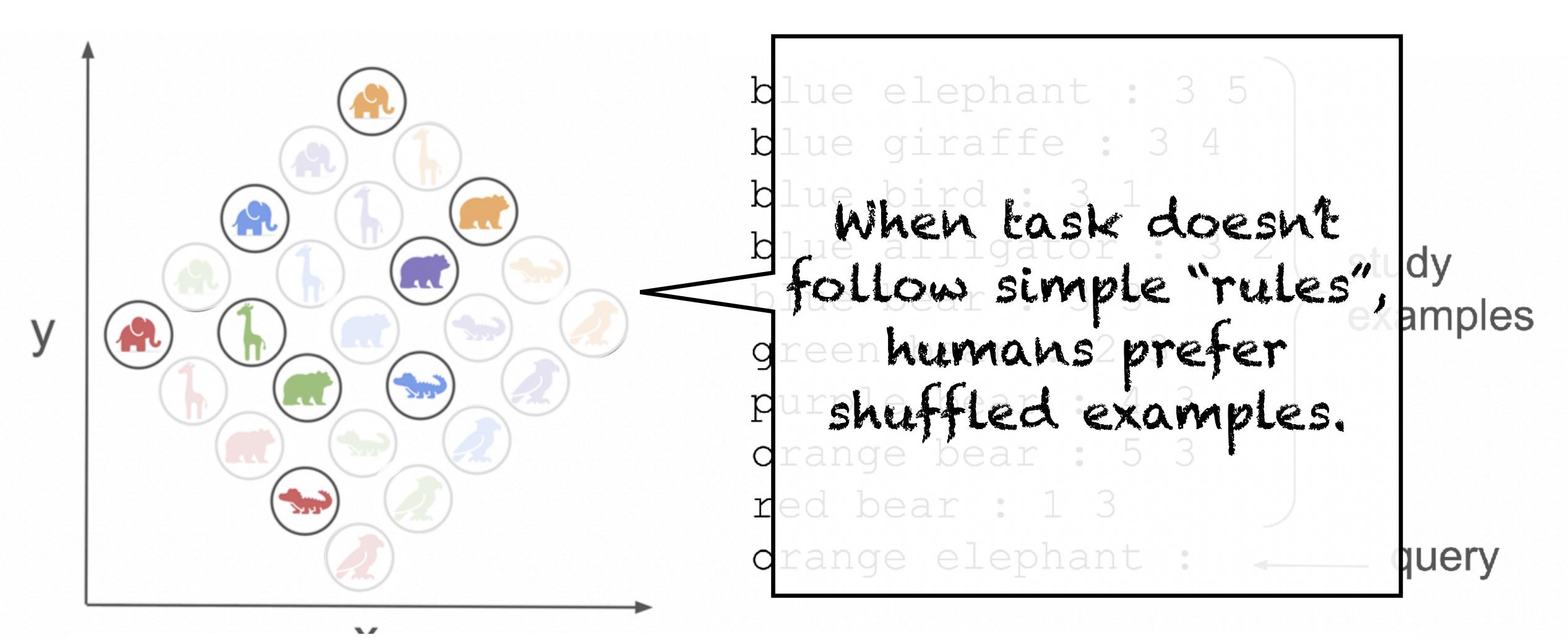
- LLMs are not merely "stochastic parrots". They often contain interpretable internal structure which makes use of modular functionally-specialized components.
- These components can sometimes resemble known brain mechanisms, e.g., mechanisms for cognitive control in the prefrontal cortex
- LLMs, off the shelf, exhibit behaviors for which we don't have good competing computational cognitive models, like the ability to learn flexibly both in context and in weights. Understanding these mechanisms and aligning them to those of humans could jointly offer new theories of human cognition and improve LLMs robustness, interpretability, and efficiency



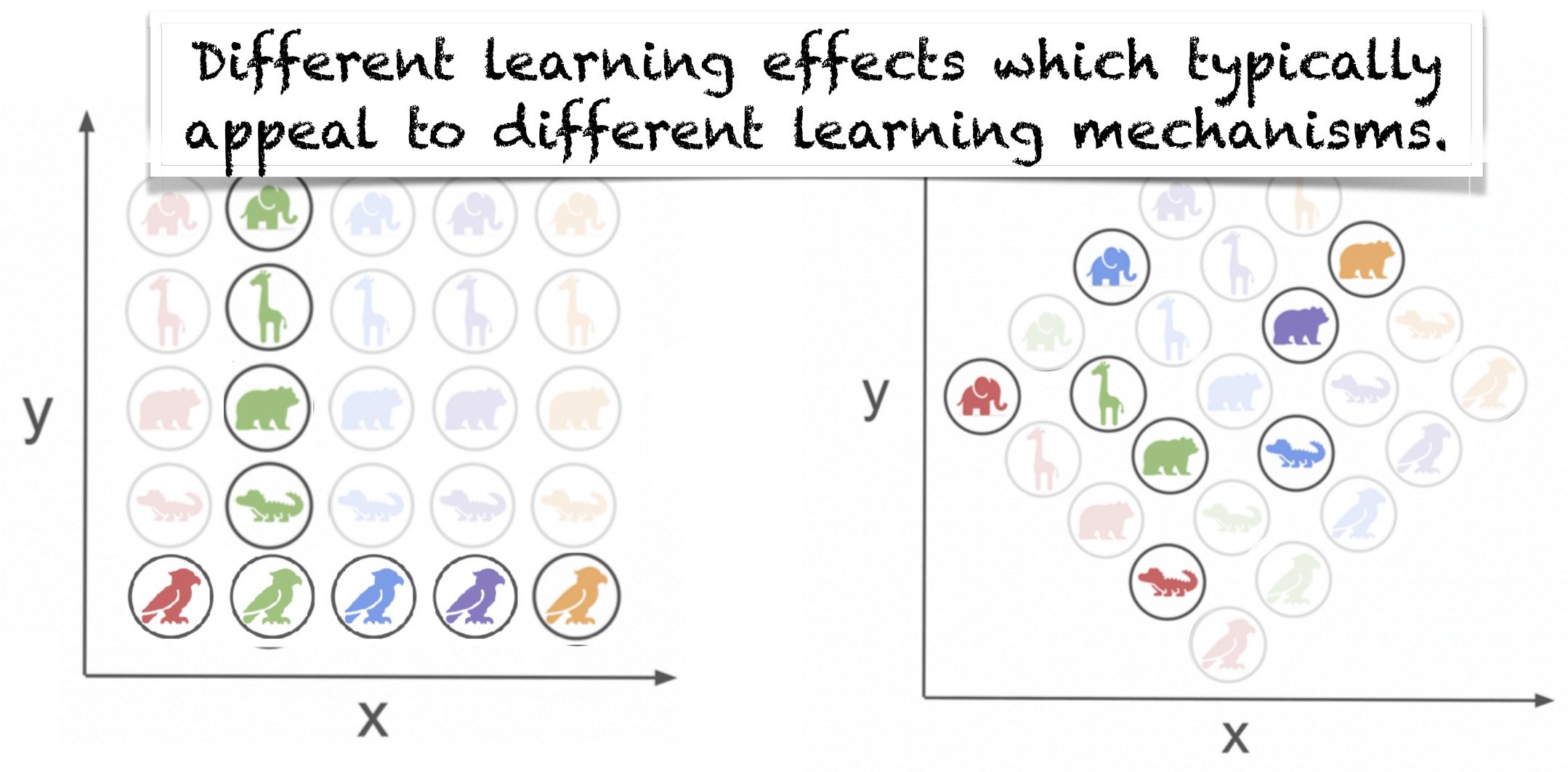
Dekker, Ronald B., Fabian Otto, and Christopher Summerfield. "Curriculum learning for human compositional generalization." Proceedings of the National Academy of Sciences 119.41 (2022): e2205582119.

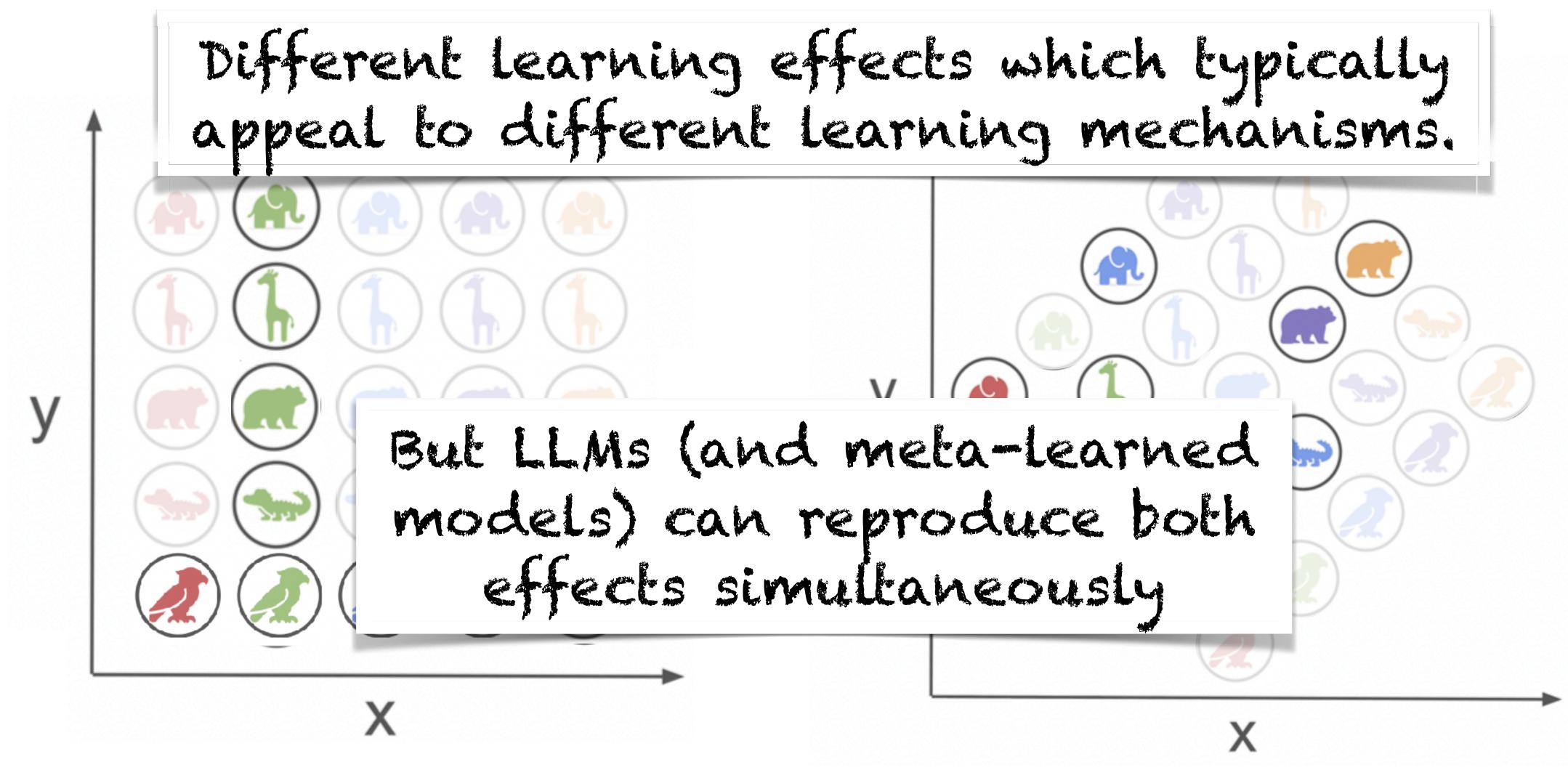


كلادر, Ronald B., Fabian Otto, and Christopher عناسات والمرادة المرادة المراد



Noh, Sharon M., et al. "Optimal sequencing during category learning: Testing a dual-learning systems perspective." Cognition 155 (2016): 23-29.





Russin, Jacob, Ellie Pavlick, and Michael J. Frank. "Human Curriculum Effects Emerge with In-Context Learning in Neural Networks." arXiv preprint arXiv:2402.08674 (2024).

Proofs of Concept

- LLMs are not merely "stochastic parrots". They often contain interpretable internal structure which makes use of modular functionally-specialized components.
- These components can sometimes resemble known brain mechanisms, e.g., mechanisms for cognitive control in the prefrontal cortex
- LLMs, off the shelf, exhibit behaviors for which we don't have good competing computational cognitive models, like the ability to learn flexibly both in context and in weights. Understanding these mechanisms and aligning them to those of humans could jointly offer new theories of human cognition and improve LLMs robustness, interpretability, and efficiency

Moving Forward

- Using LLMs as cognitive models has the potential to unlock a virtuous cycle that both improves our understanding of human cognition and improves Al
- Harnessing this requires:
 - Open source models for research which match the scale and performance of proprietary models
 - Investment in research with a long time horizon, which falls outside of commercial priorities (decades, not years)
 - (Even more) open-minded, interdisciplinary collaboration

Thank you!