

Neuro-vector-symbolic architectures: Toward Computationally Efficient Machine Learning and Reasoning utilizing In-Memory Computing

Abbas Rahimi

IBM Research-Zurich

Mondays in Memory (MiM) Webinar

Feb. 3, 2025

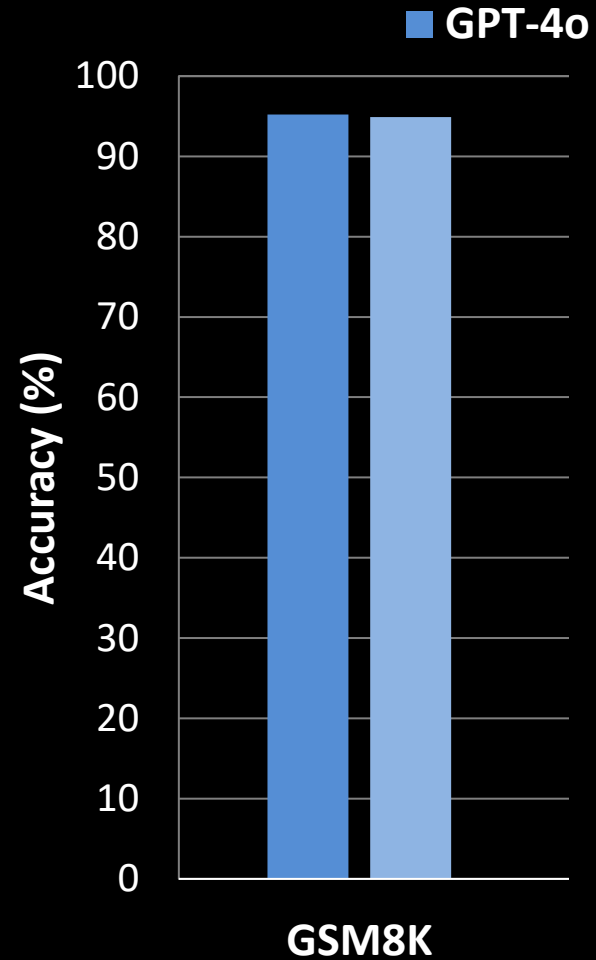
AI today



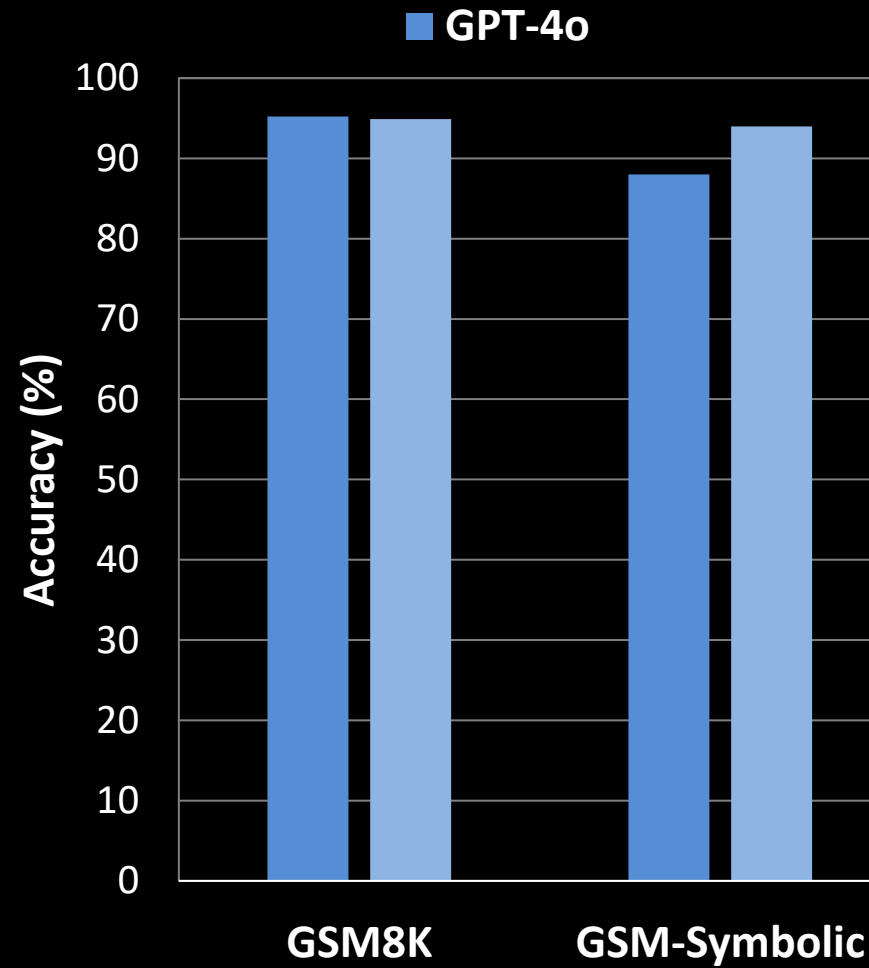
Struggling with out-of-distribution (OOD)



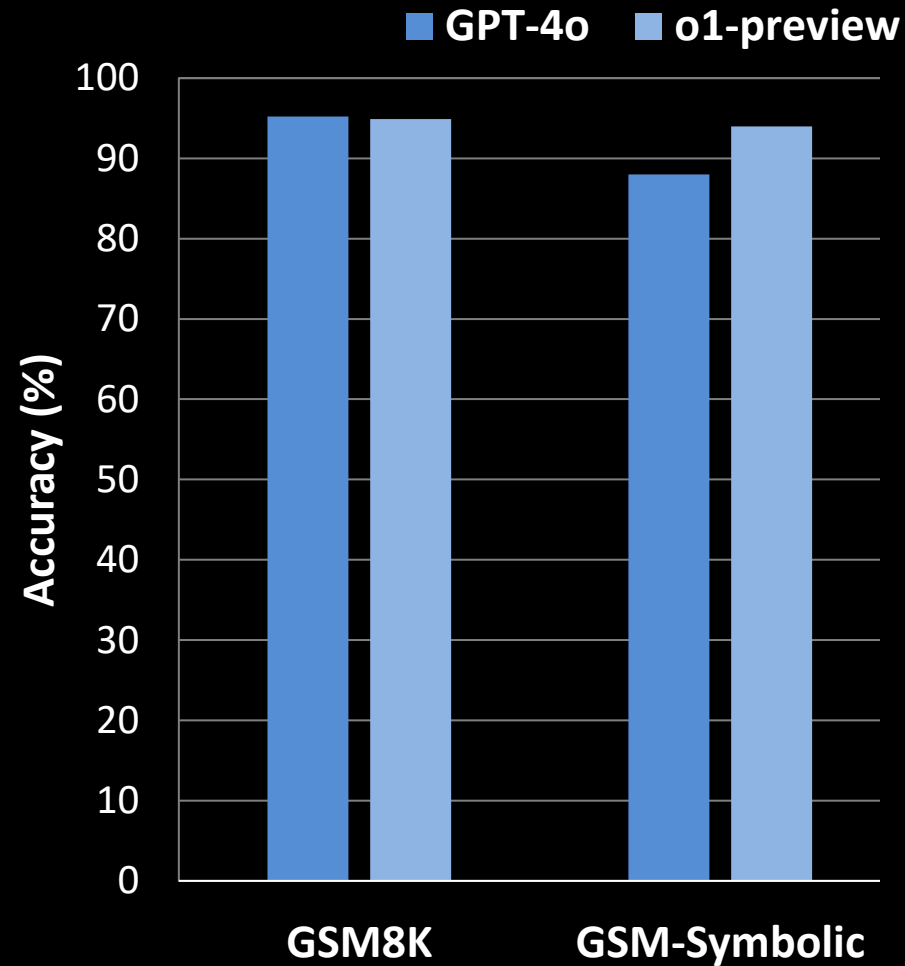
Struggling with out-of-distribution (OOD)



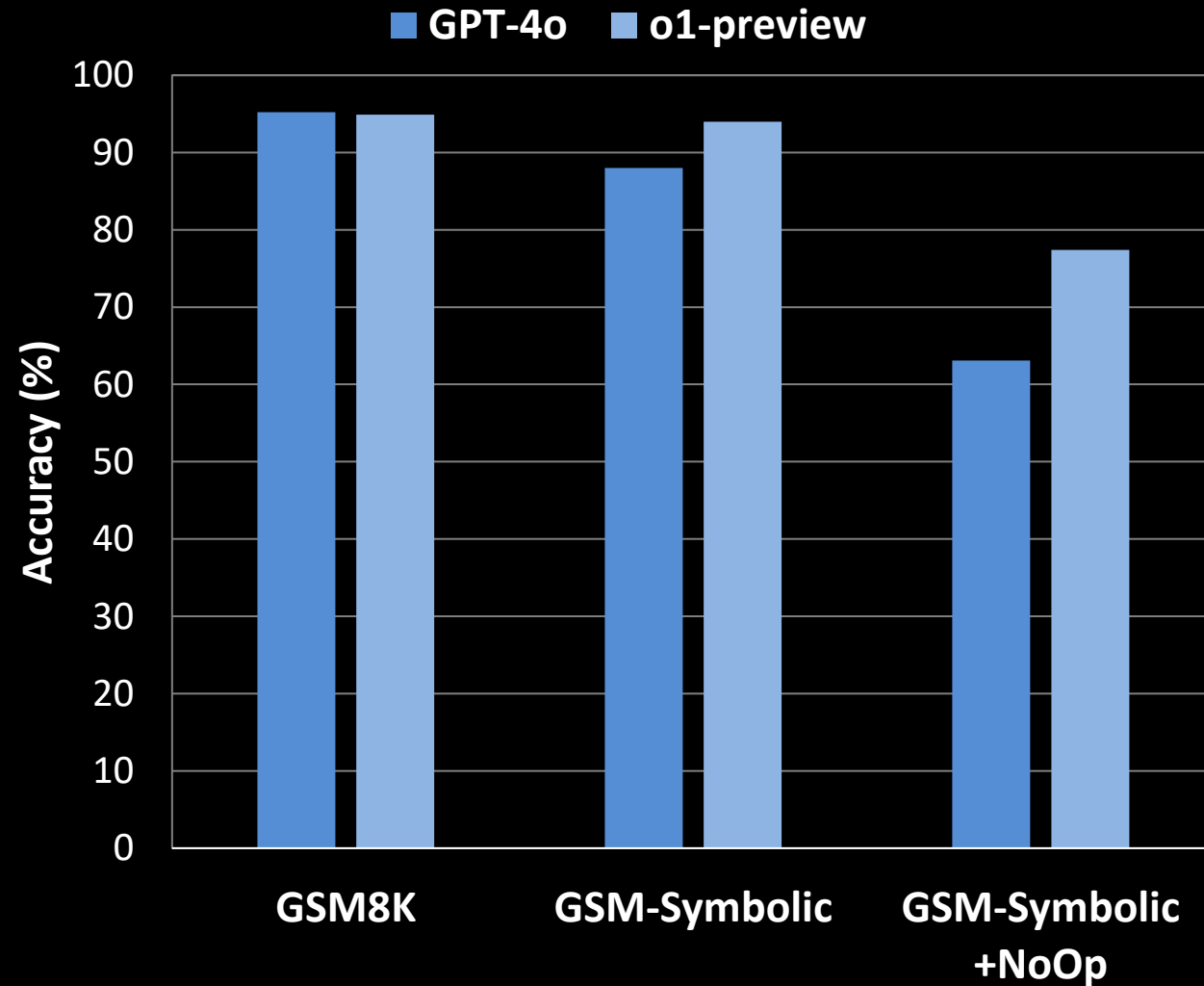
Struggling with out-of-distribution (OOD)



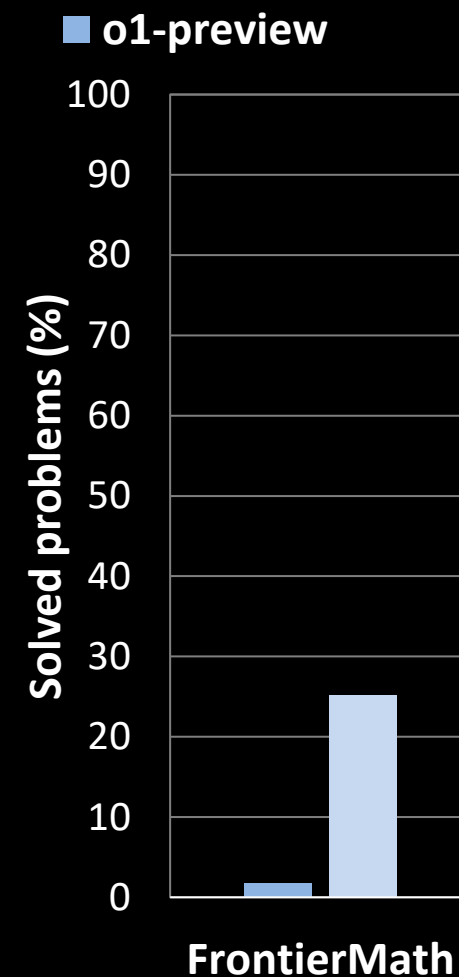
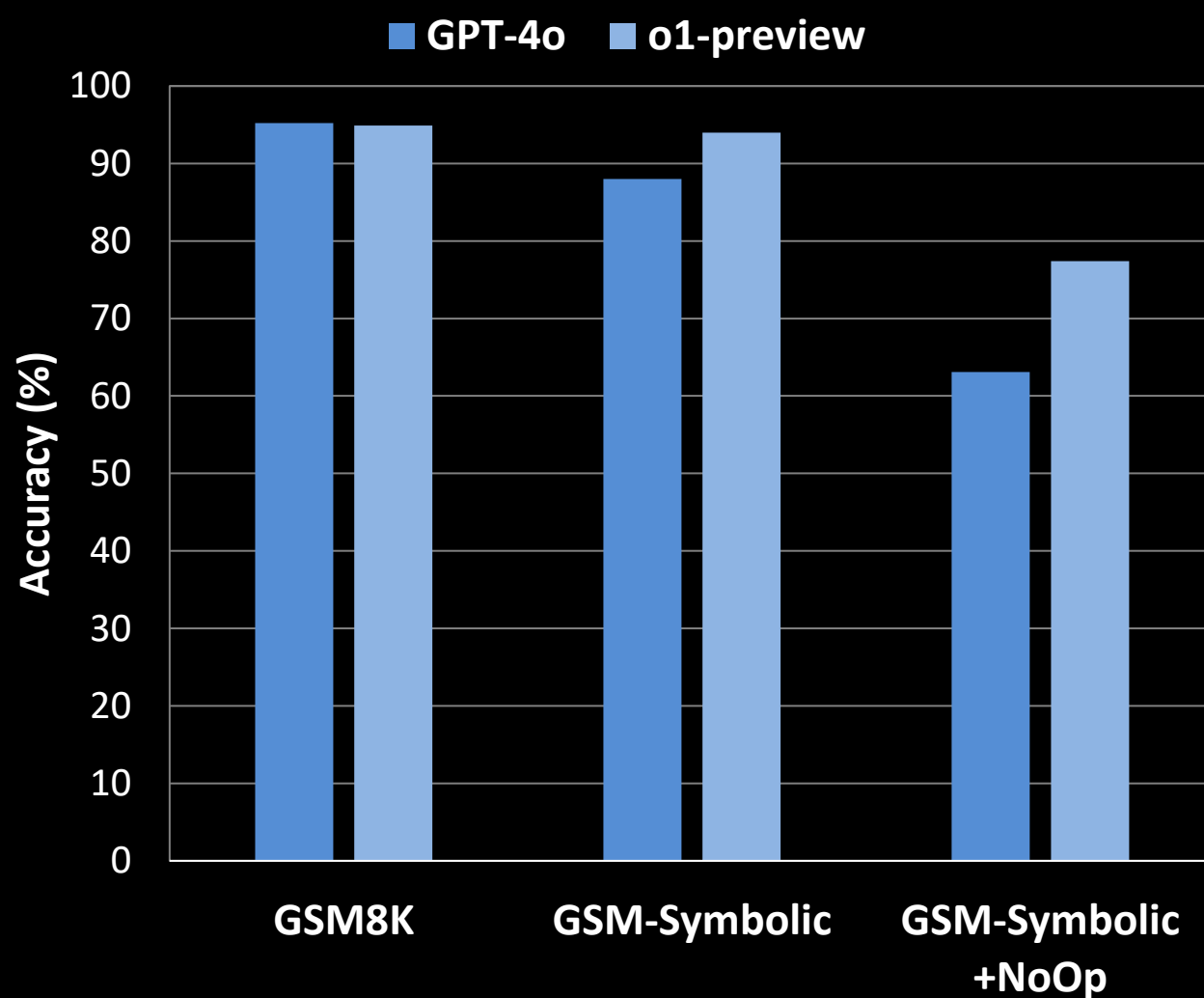
Struggling with out-of-distribution (OOD)



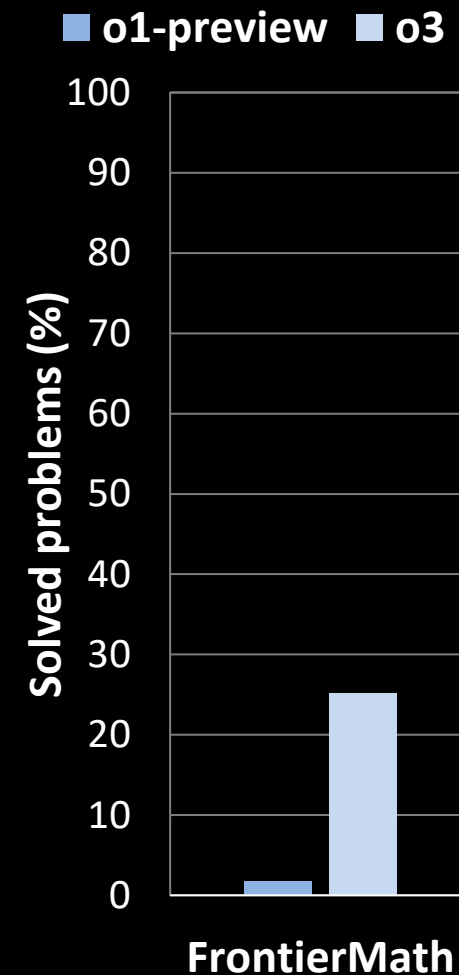
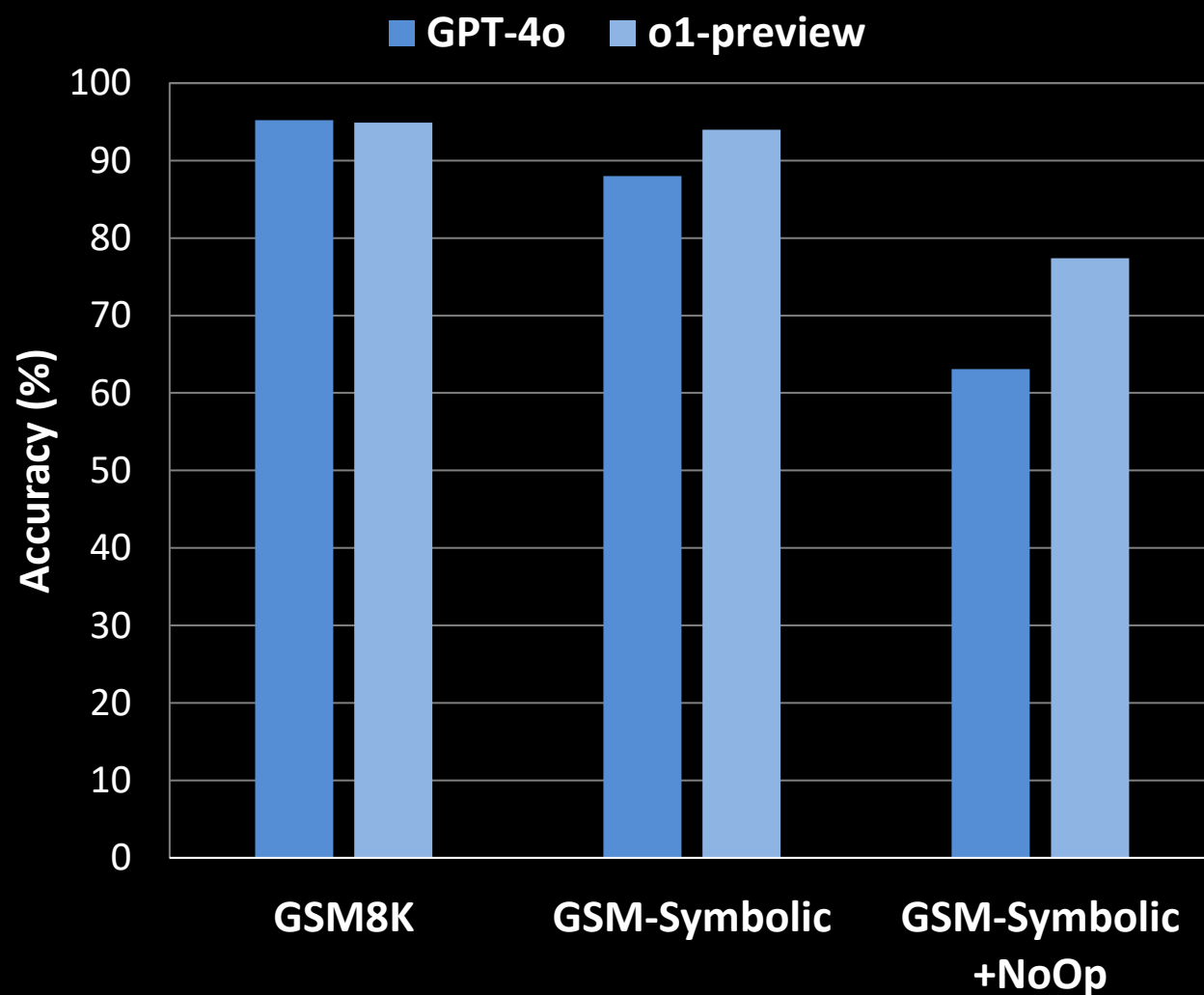
Struggling with out-of-distribution (OOD)



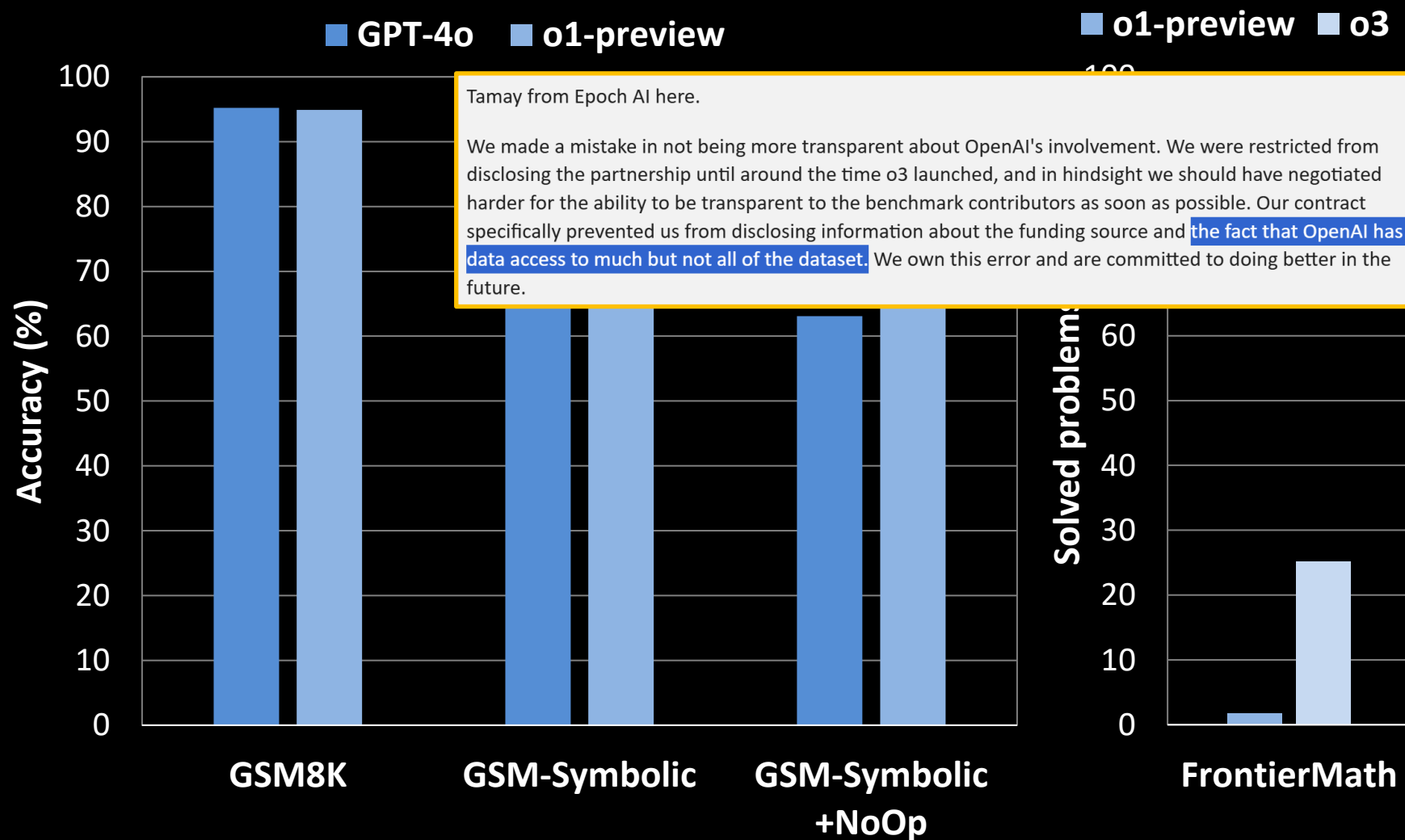
Struggling with out-of-distribution (OOD)



Struggling with out-of-distribution (OOD)



Struggling with out-of-distribution (OOD)



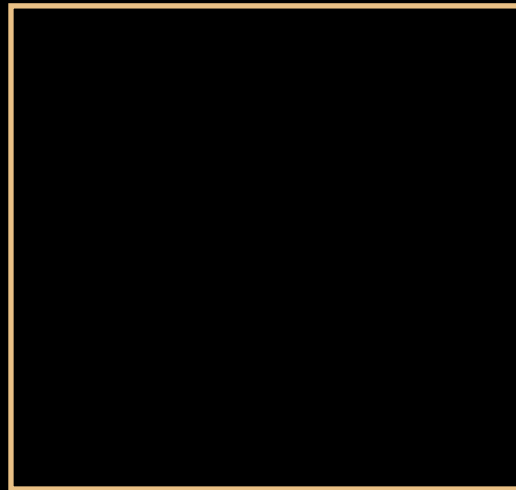
Even for in-distribution data, Transformers are **sample inefficient learners** for solving compositions

$$\begin{array}{r} \times \quad 8 \ 0 \ 4 \ 9 \\ \quad \quad \quad 1 \ 1 \\ \hline \quad \quad 8 \ 0 \ 4 \ 9 \\ 8 \ 0 \ 4 \ 9 \\ \hline 8 \ 8 \ 5 \ 3 \ 9 \end{array}$$

Even for in-distribution data, Transformers are **sample inefficient learners** for solving compositions

$$\begin{array}{r} \times \quad 8 \ 0 \ 4 \ 9 \\ \quad \quad \quad 1 \ 1 \\ \hline \quad \quad 8 \ 0 \ 4 \ 9 \\ 8 \ 0 \ 4 \ 9 \\ \hline 8 \ 8 \ 5 \ 3 \ 9 \end{array}$$

Step-by-step
Multiplication



Even for in-distribution data, Transformers are **sample inefficient learners** for solving compositions

$$\begin{array}{r} \times \quad \boxed{8 \ 0 \ 4 \ 9} \\ \quad \quad \boxed{1 \ 1} \\ \hline \quad \boxed{8 \ 0 \ 4 \ 9} \\ 8 \ 0 \ 4 \ 9 \\ \hline 8 \ 8 \ 5 \ 3 \ 9 \end{array}$$

Step-by-step
Multiplication

Digit Multipli.

Even for in-distribution data, Transformers are **sample inefficient learners** for solving compositions

$$\begin{array}{r} \times \quad \boxed{8049} \\ \quad \quad \boxed{1}1 \\ \hline \quad 8049 \\ \boxed{8049} \\ \hline 88539 \end{array}$$

Step-by-step
Multiplication

Digit Multipli.

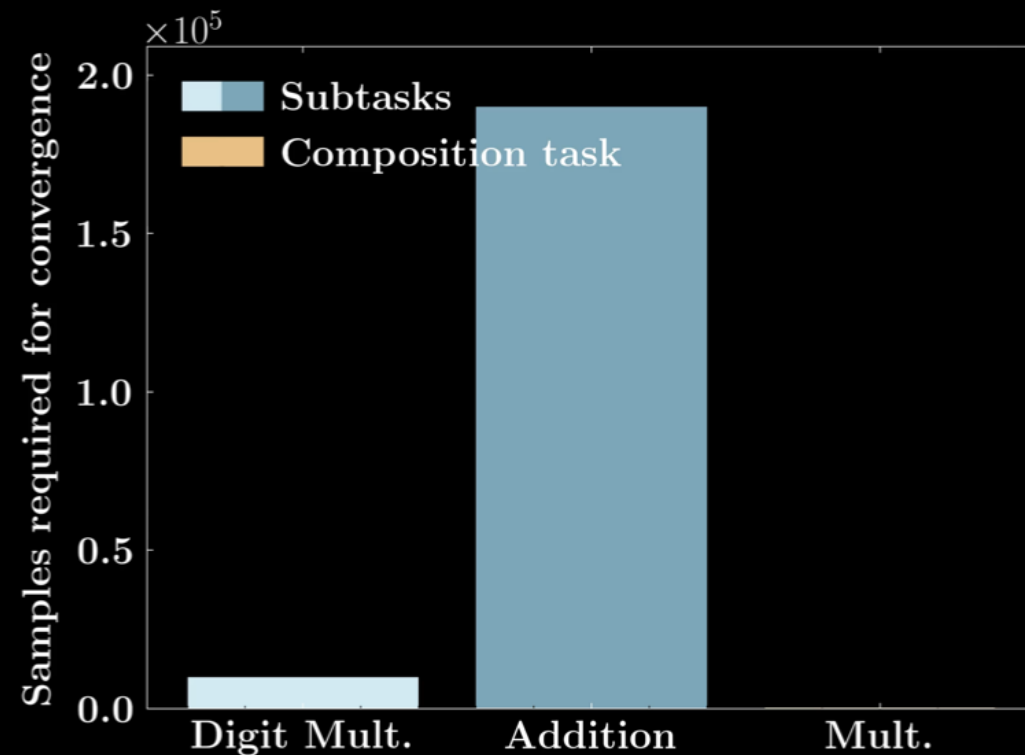
Digit Multipli.

Even for in-distribution data, Transformers are **sample inefficient learners** for solving compositions

$$\begin{array}{r} \times \quad 8 \ 0 \ 4 \ 9 \\ \quad \quad 1 \ 1 \\ \hline \boxed{\begin{array}{r} 8 \ 0 \ 4 \ 9 \\ 8 \ 0 \ 4 \ 9 \end{array}} \\ \hline 8 \ 8 \ 5 \ 3 \ 9 \end{array}$$

Step-by-step
Multiplication

Digit Multipli.
Digit Multipli.
Addition



Even for in-distribution data, Transformers are **sample inefficient learners** for solving compositions

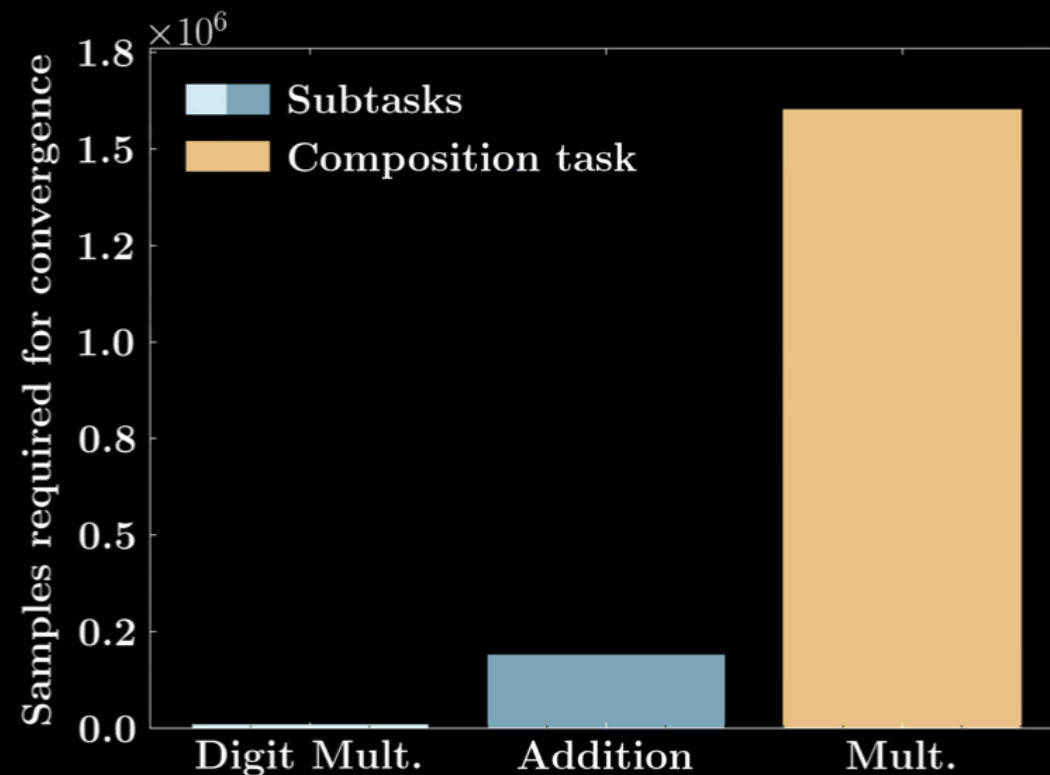
$$\begin{array}{r} \times \quad 8049 \\ \quad 11 \\ \hline 8049 \\ 8049 \\ \hline \boxed{88539} \end{array}$$

Step-by-step
Multiplication

Digit Multipli.

Digit Multipli.

Addition

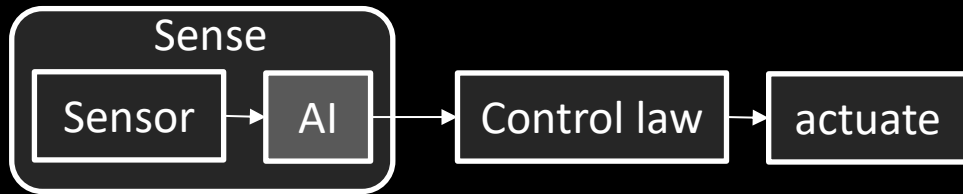


Edward Lee: “Certainty or Intelligence: Pick One!”

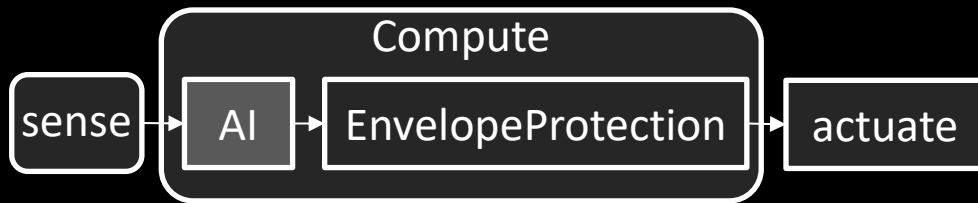
How to deal with contradictory goals?

Lessons learned from cyber-physical systems:

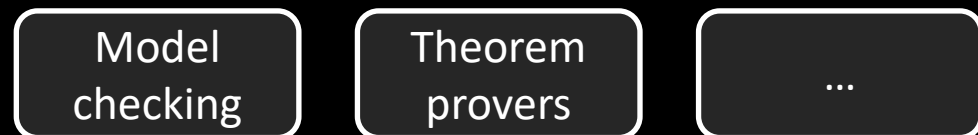
1) AI as sensor



2) Envelope protection



3) Rely on formal methods

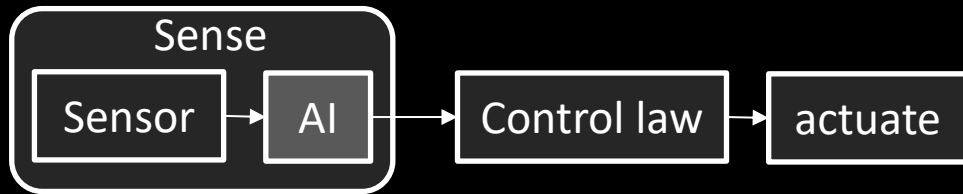


Edward Lee: “Certainty or Intelligence: Pick One!”

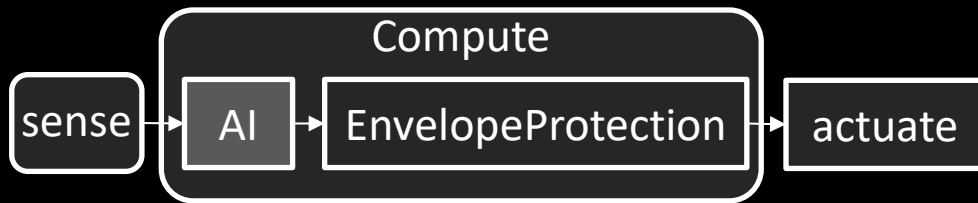
How to deal with contradictory goals?

Lessons learned from cyber-physical systems:

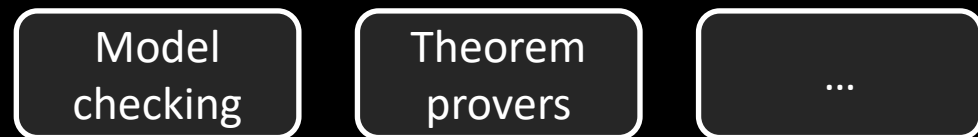
1) AI as sensor



2) Envelope protection



3) Rely on formal methods



Neuro-symbolic (NeSy) AI:

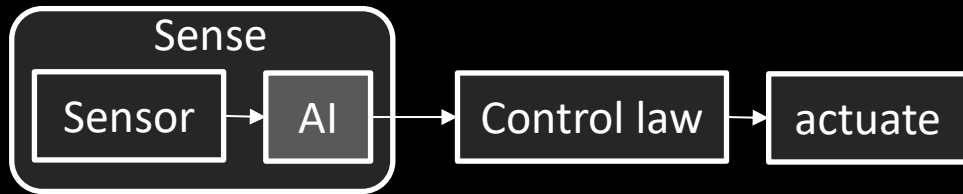
1. Symbolic Neuro Symbolic
2. Symbolic[Neuro]
3. Neuro;Symbolic
4. Neuro:Symbolic→Neuro
5. Neuro_{Symbolic}
6. Neuro[Symbolic]

Edward Lee: “Certainty or Intelligence: Pick One!”

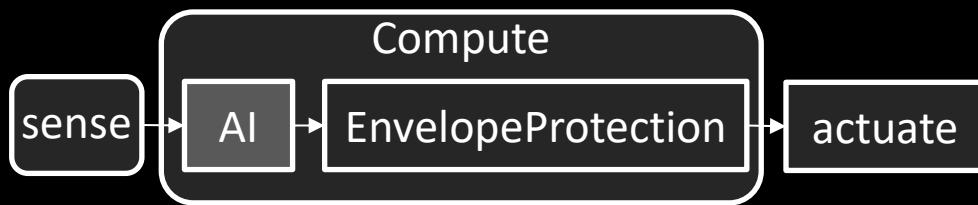
How to deal with contradictory goals?

Lessons learned from cyber-physical systems:

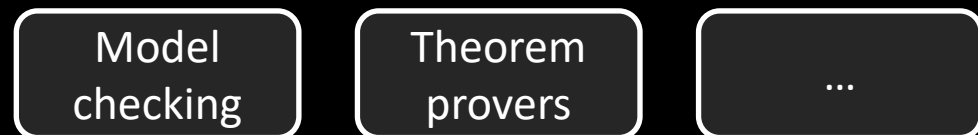
1) AI as sensor



2) Envelope protection



3) Rely on formal methods



Neuro-symbolic (NeSy) AI:

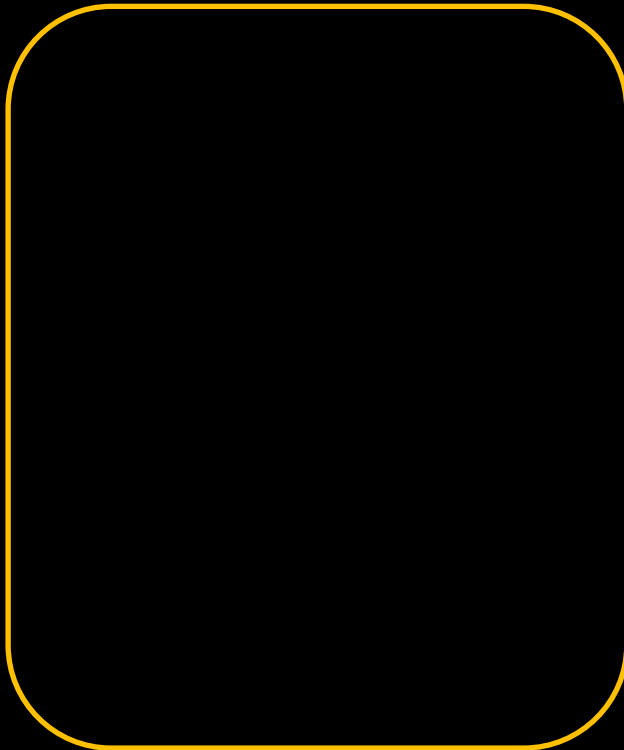
1. Symbolic Neuro Symbolic
2. Symbolic[Neuro]
3. Neuro;Symbolic
4. Neuro:Symbolic→Neuro
5. Neuro_{Symbolic}
6. Neuro[Symbolic]

Neuro-vector-symbolic architectures (NVSA)

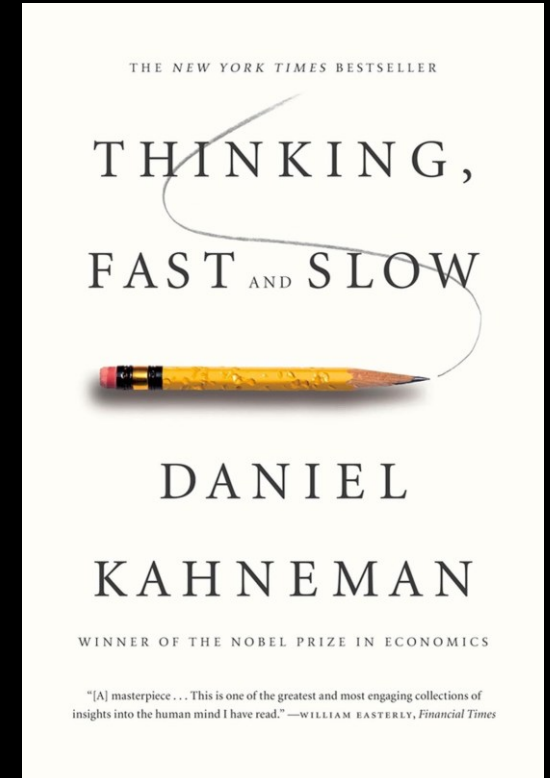
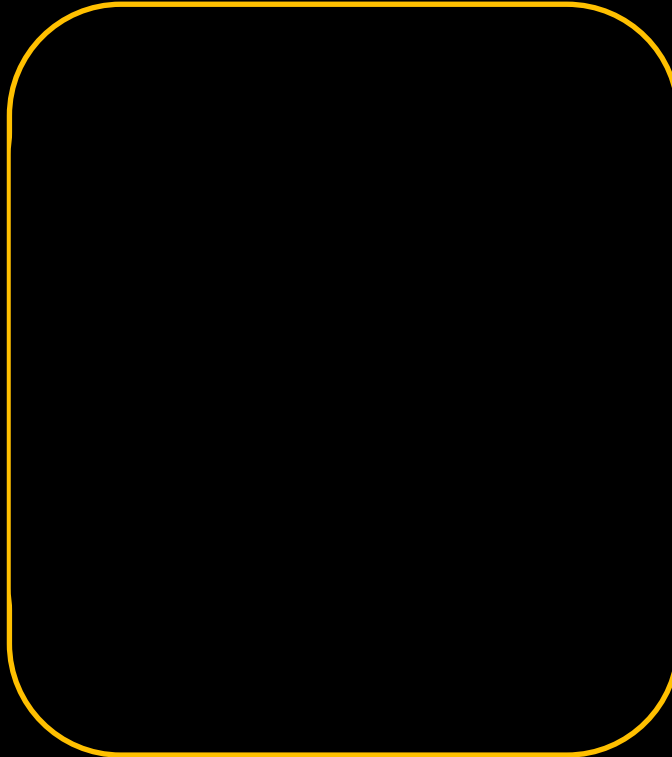
NVSA complements neural nets with vector-symbolic architectures in a unified framework

NVSA complements neural nets with vector-symbolic architectures in a unified framework

System 1: Perception



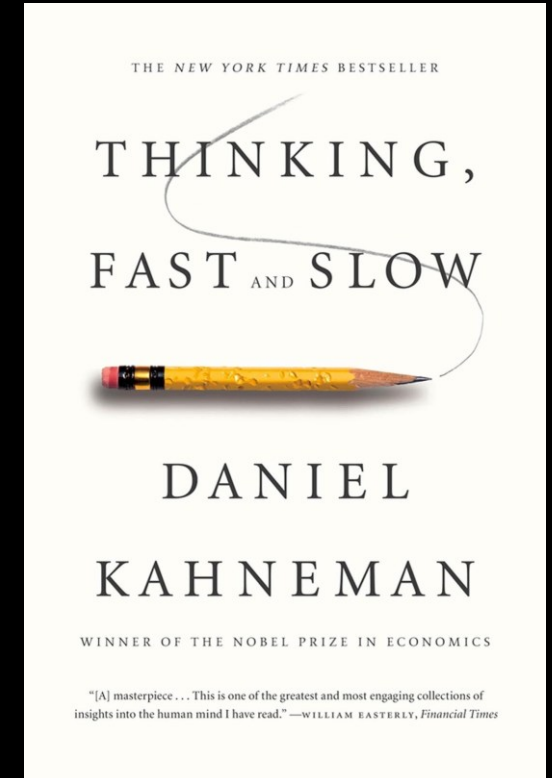
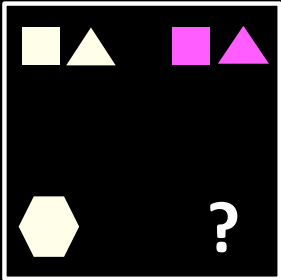
System 2: Reasoning



NVSA complements neural nets with vector-symbolic architectures in a unified framework

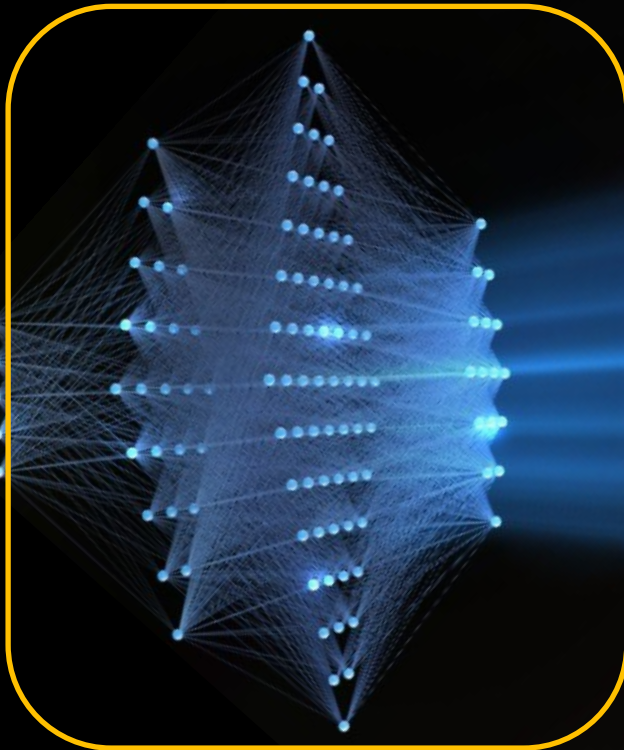
System 1: Perception

System 2: Reasoning

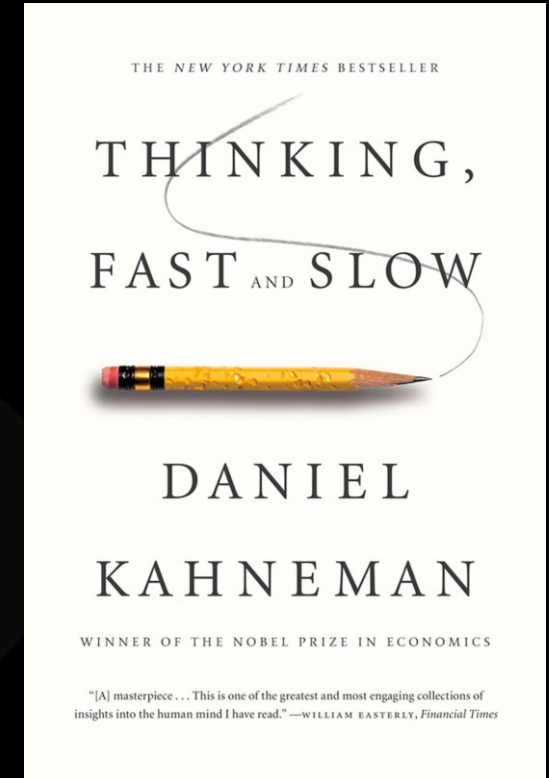
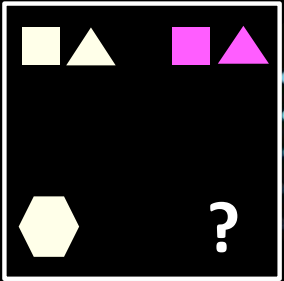


NVSA complements neural nets with vector-symbolic architectures in a unified framework

System 1: Perception

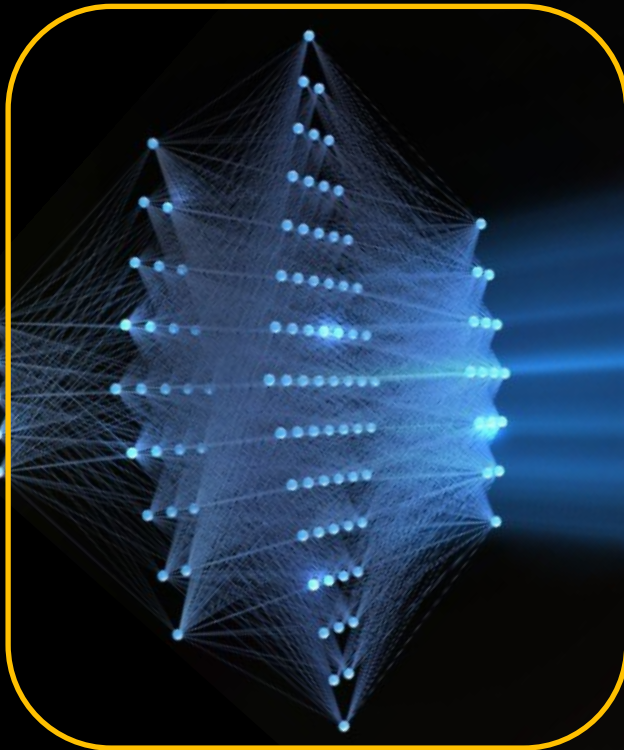


System 2: Reasoning

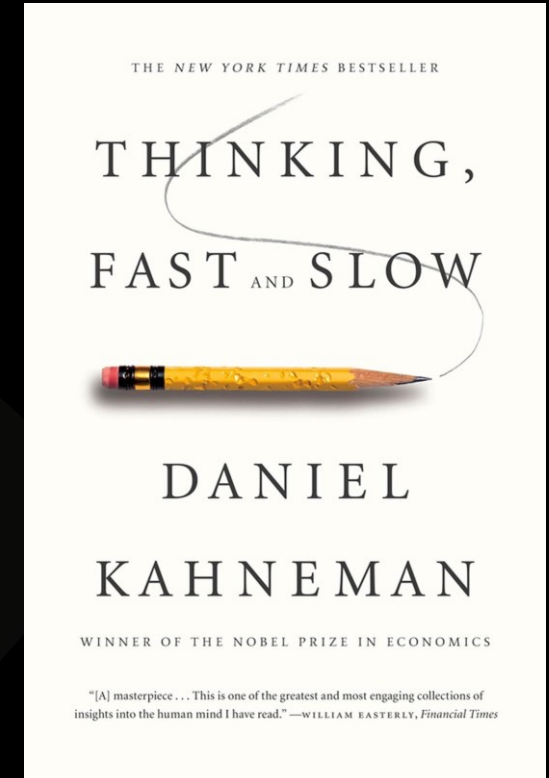
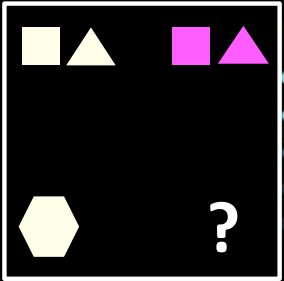
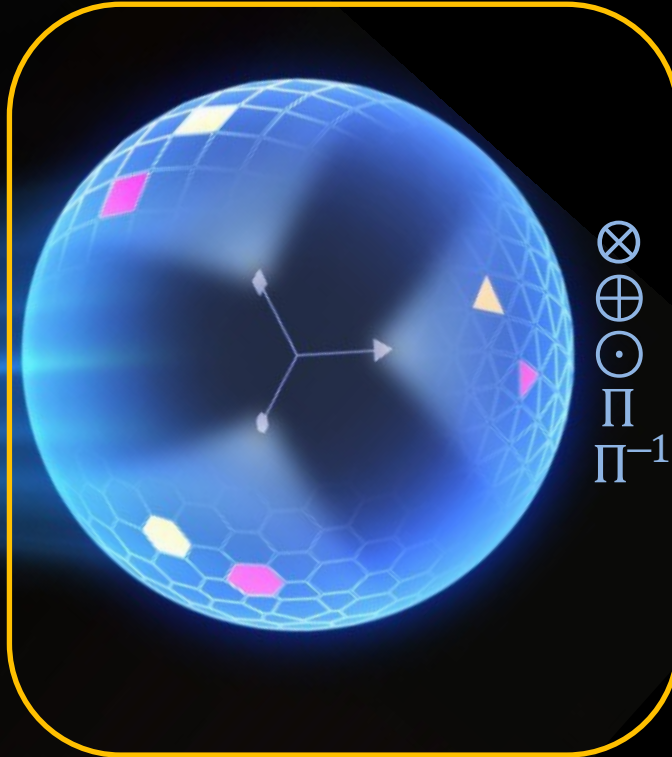


NVSA complements neural nets with vector-symbolic architectures in a unified framework

System 1: Perception

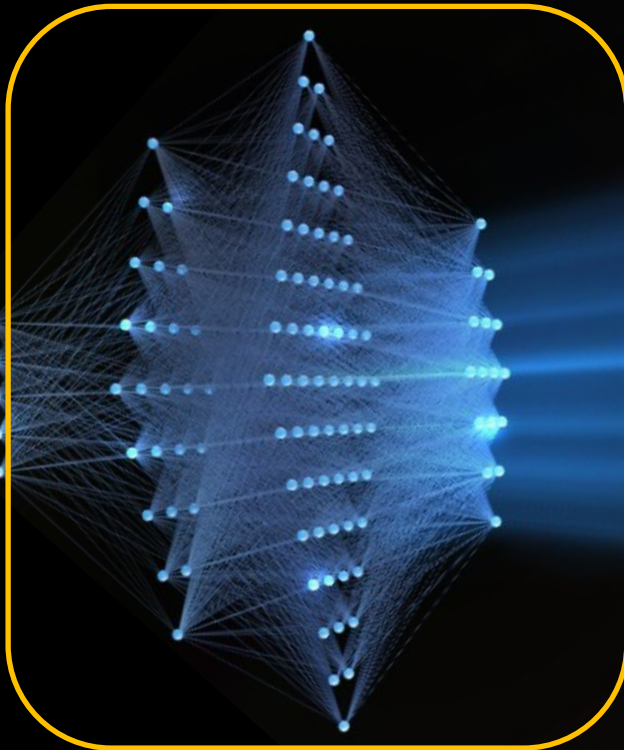


System 2: Reasoning

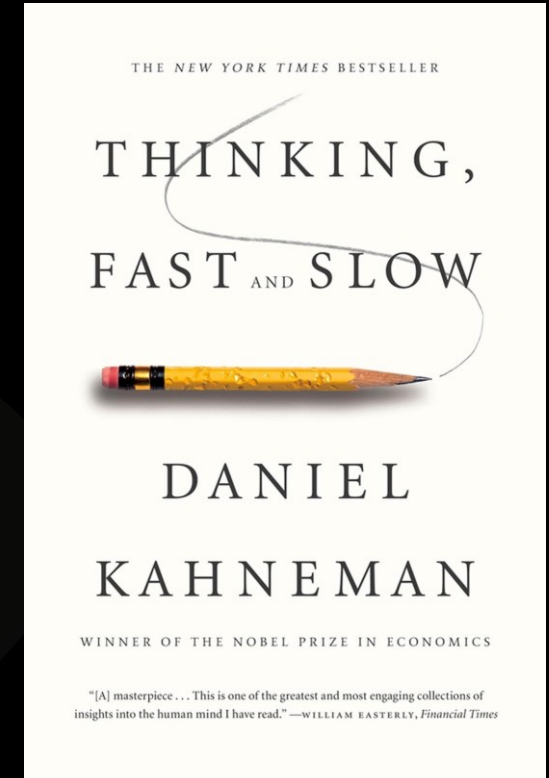
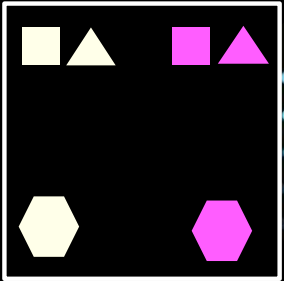
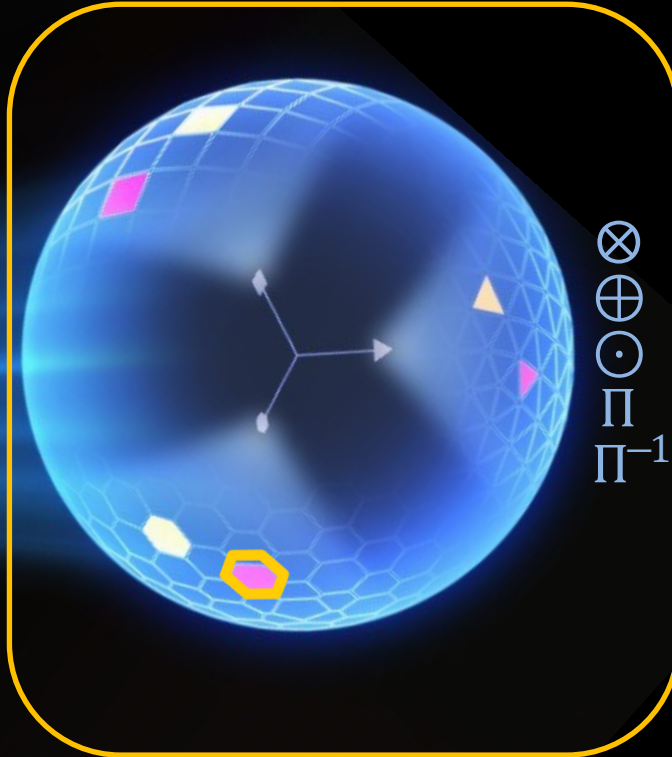


NVSA complements neural nets with vector-symbolic architectures in a unified framework

System 1: Perception

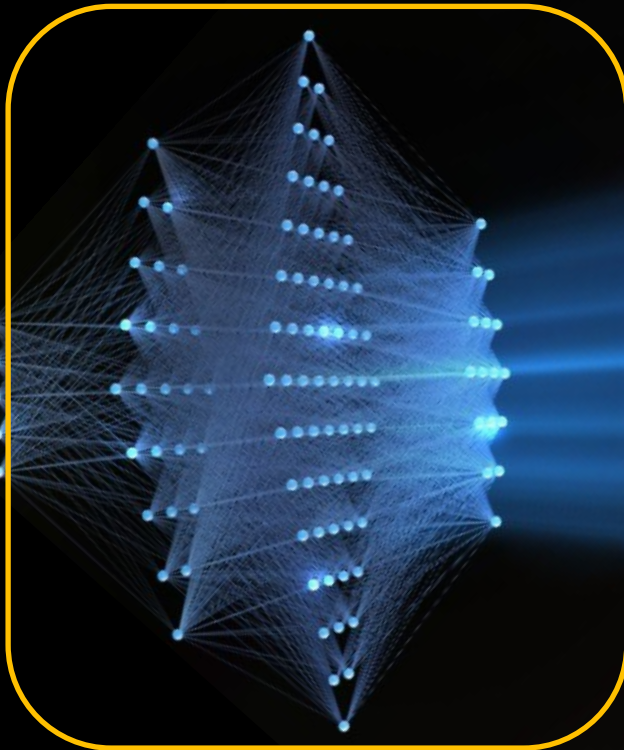


System 2: Reasoning

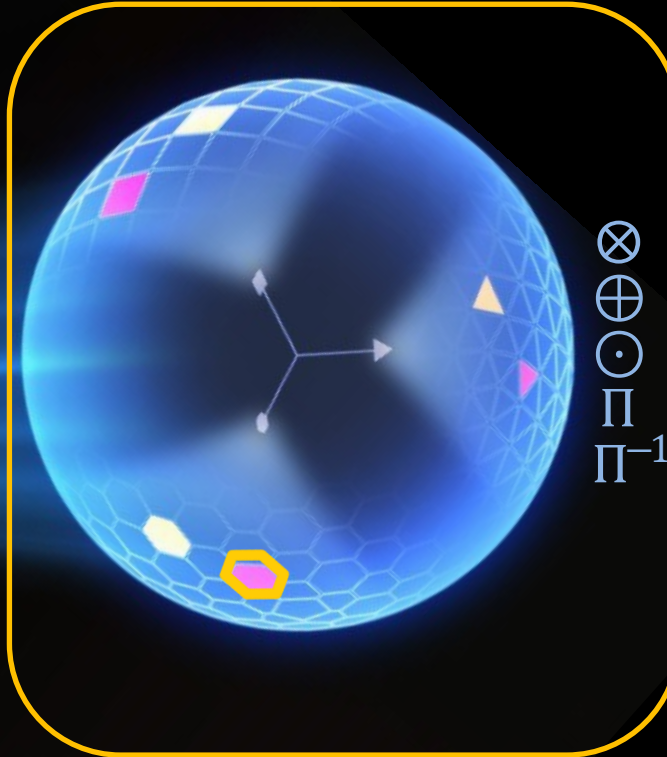


NVSA complements neural nets with vector-symbolic architectures in a unified framework

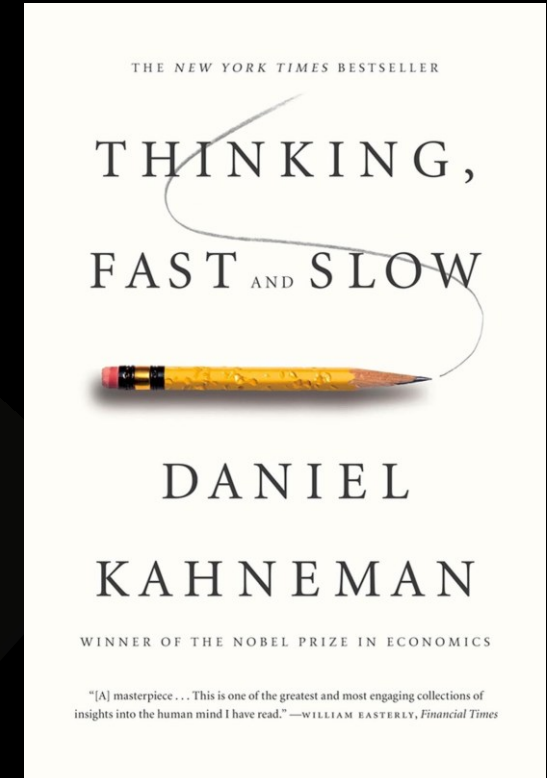
System 1: Perception



System 2: Reasoning

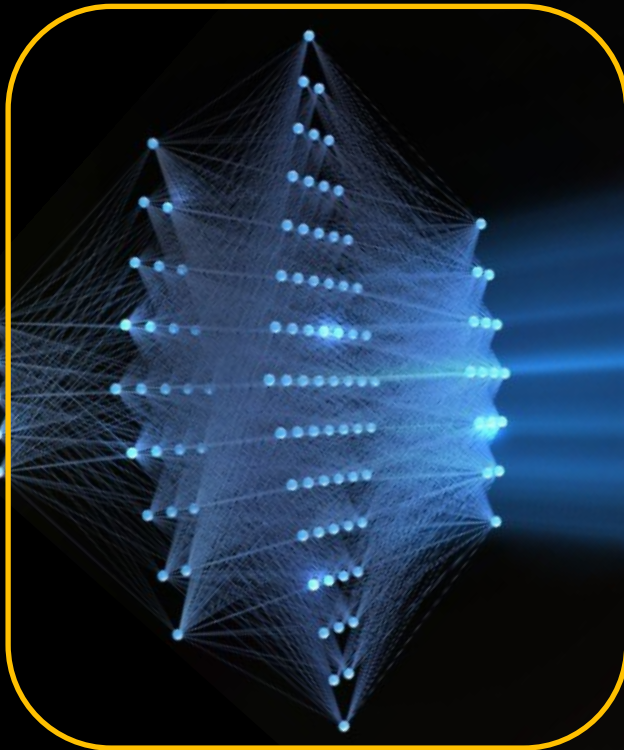


Goal: promoting compositionally structured representation in service of reasoning



NVSA complements neural nets with vector-symbolic architectures in a unified framework

System 1: Perception

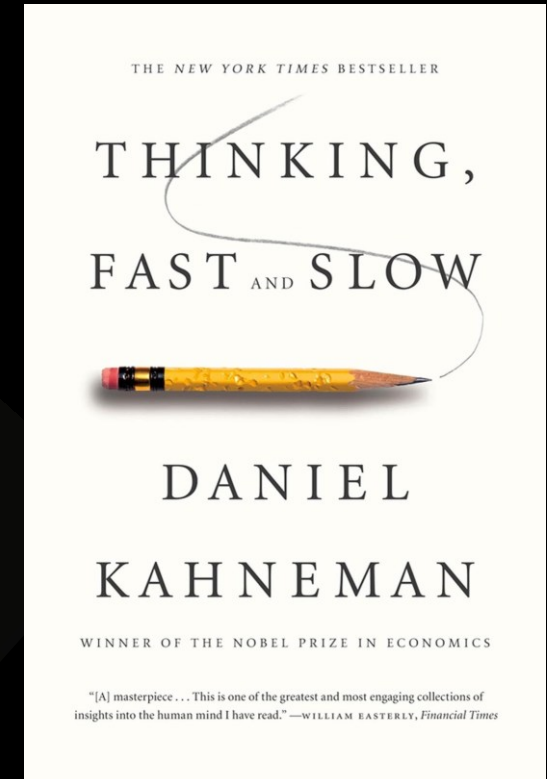
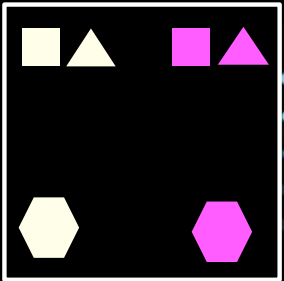


Goal: promoting compositionally structured representation in service of reasoning

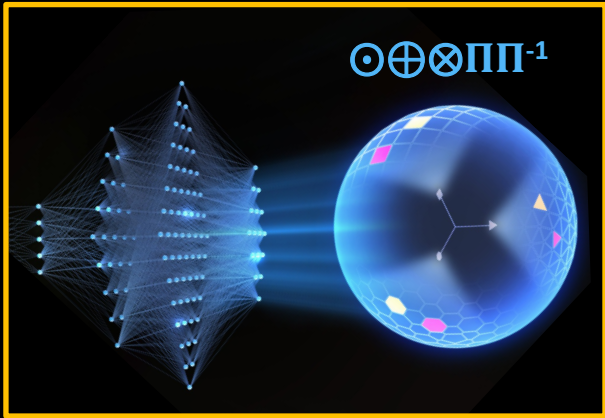
System 2: Reasoning



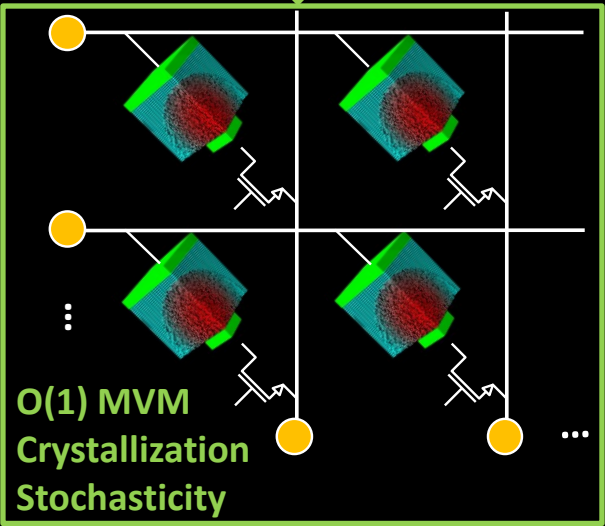
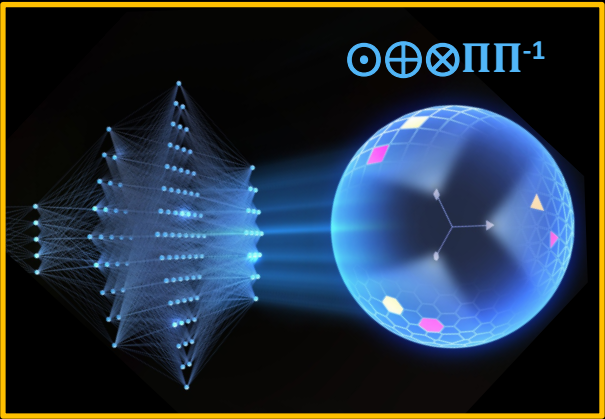
Goal: Enabling reliable yet scalable reasoning via distributed representations



**NVSA: From robust and efficient
inference to learning and reasoning**

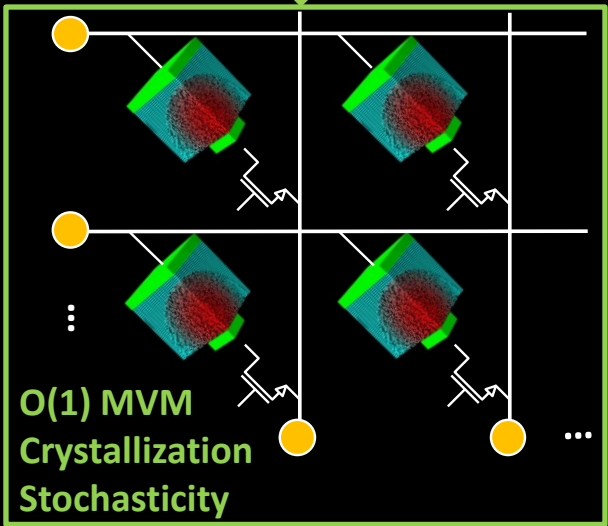
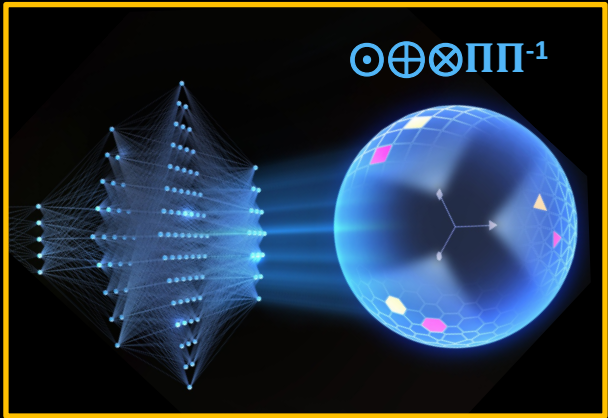


NVSA: From robust and efficient inference to learning and reasoning

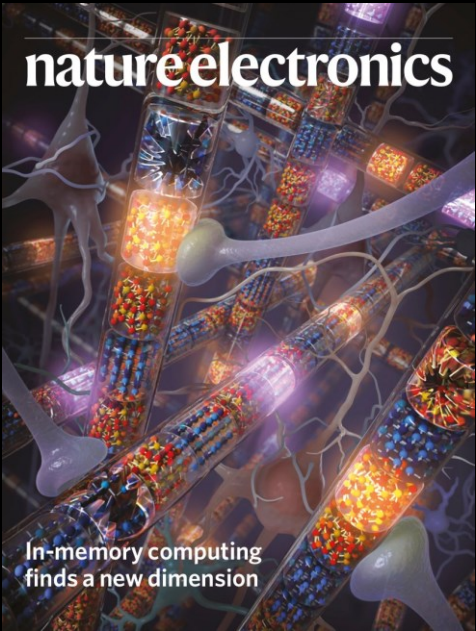


HW/SW realization is benefited
from in-memory computing

NVSA: From robust and efficient inference to learning and reasoning

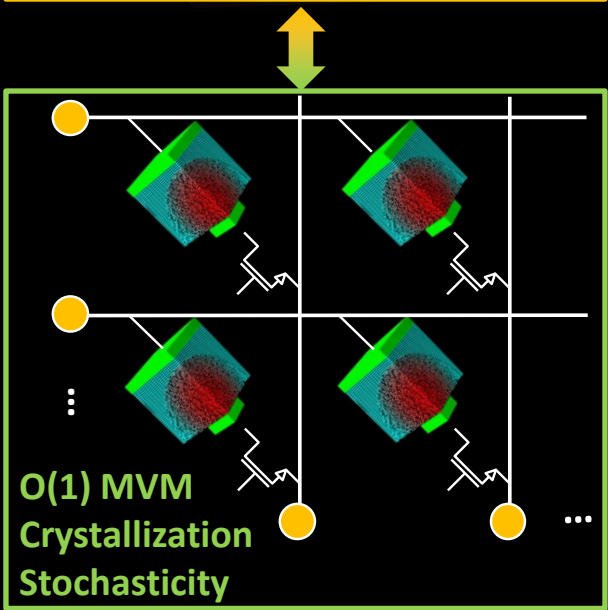
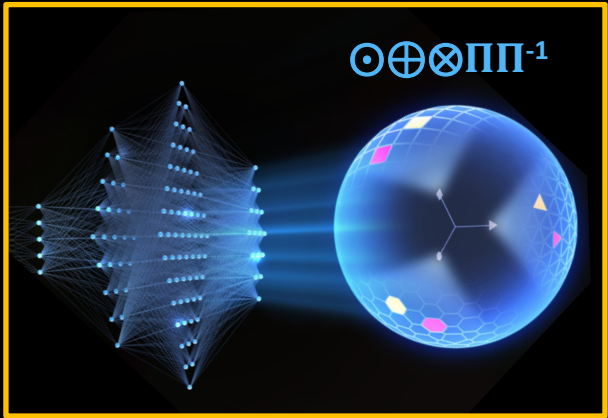


HW/SW realization is benefited from in-memory computing

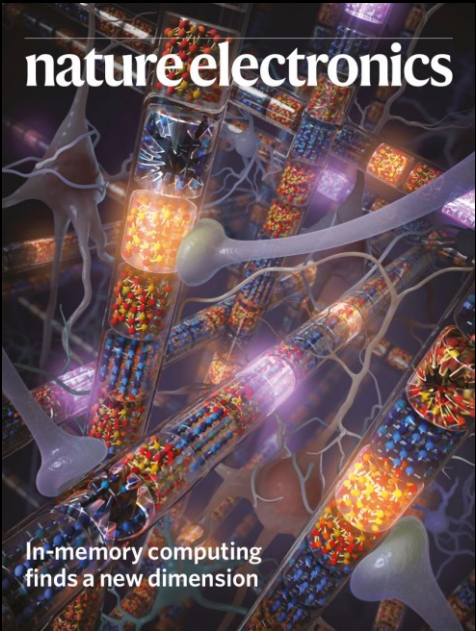


Robust VSA inference

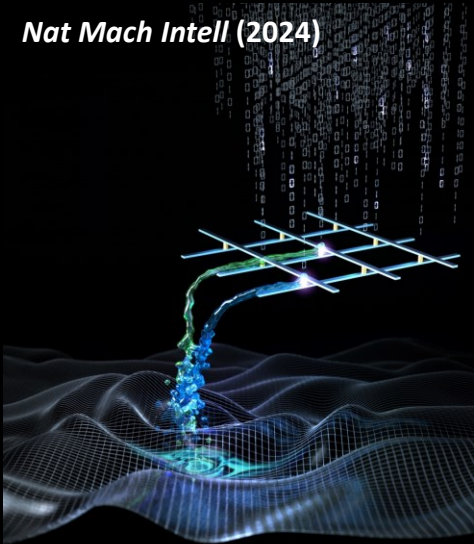
NVSA: From robust and efficient inference to learning and reasoning



HW/SW realization is benefited from in-memory computing

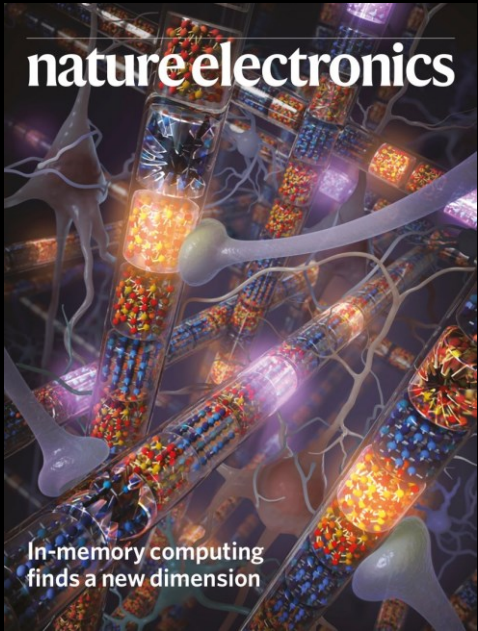
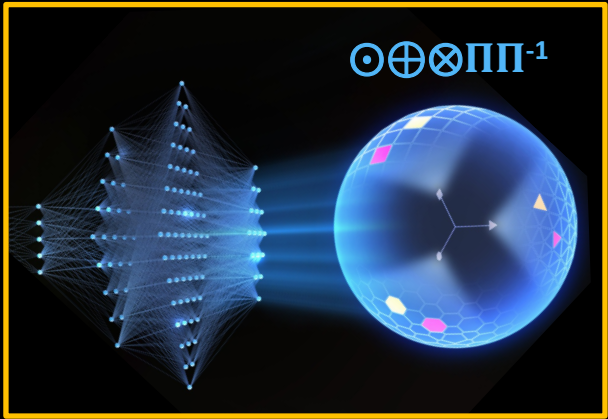


Robust VSA inference

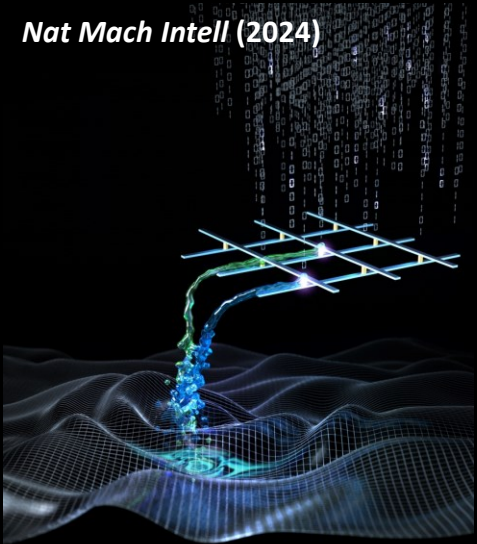


In-memory kernel approximation

NVSA: From robust and efficient inference to learning and reasoning



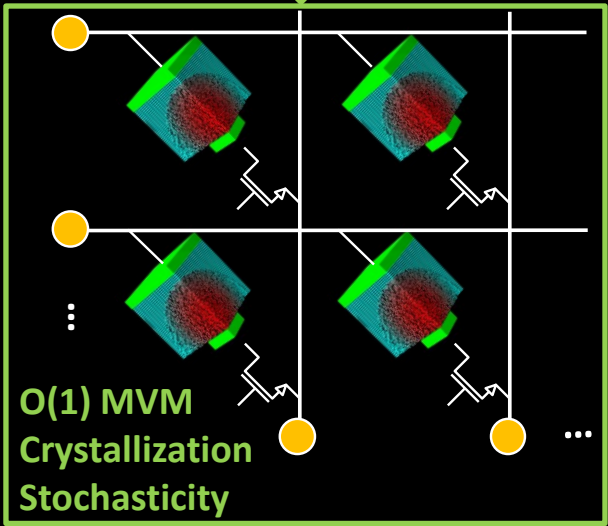
Robust VSA inference



In-memory kernel approximation

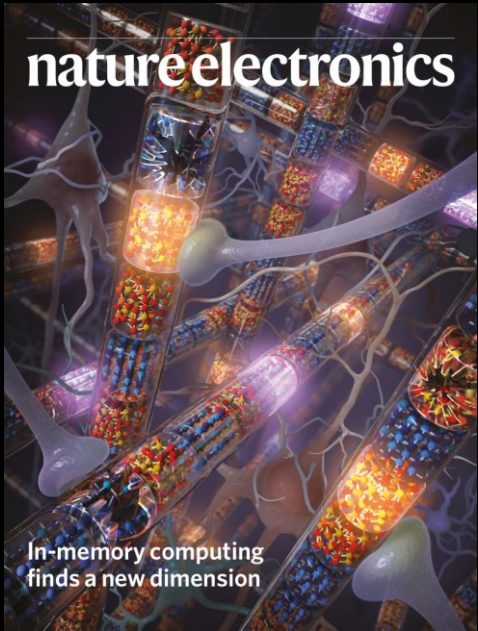
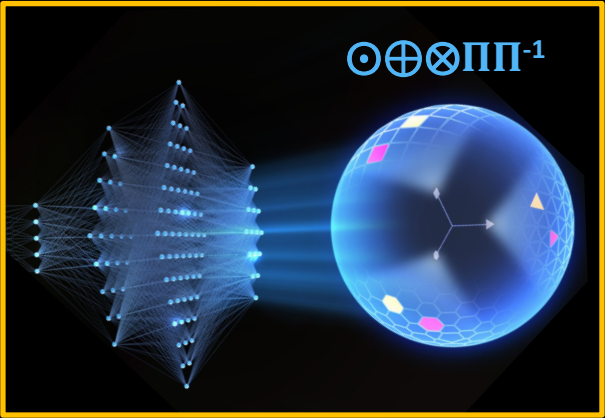


3D stack meets Mixture of Experts

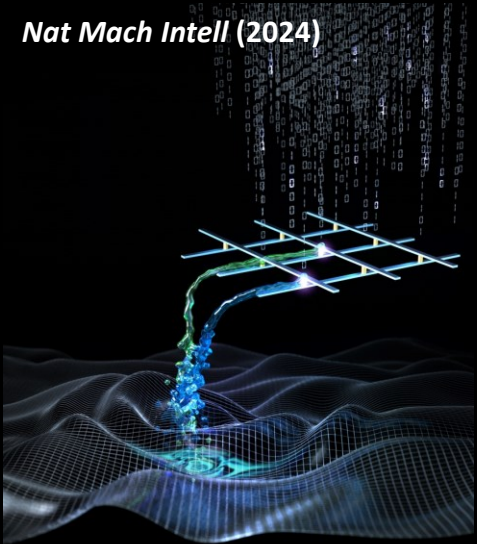


HW/SW realization is benefited from in-memory computing

NVSA: From robust and efficient inference to learning and reasoning



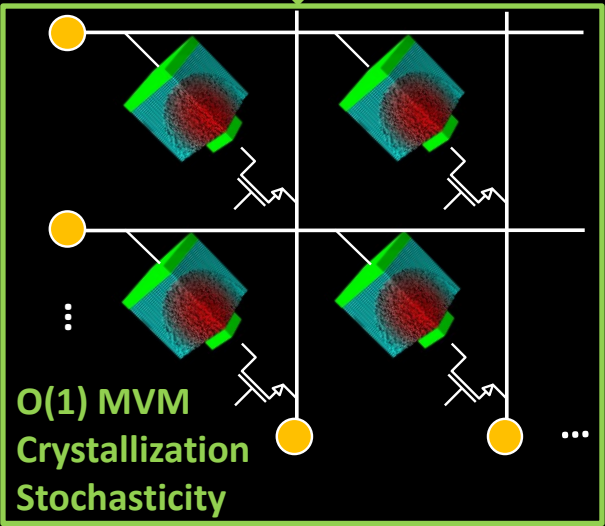
Robust VSA inference



In-memory kernel approximation



3D stack meets Mixture of Experts

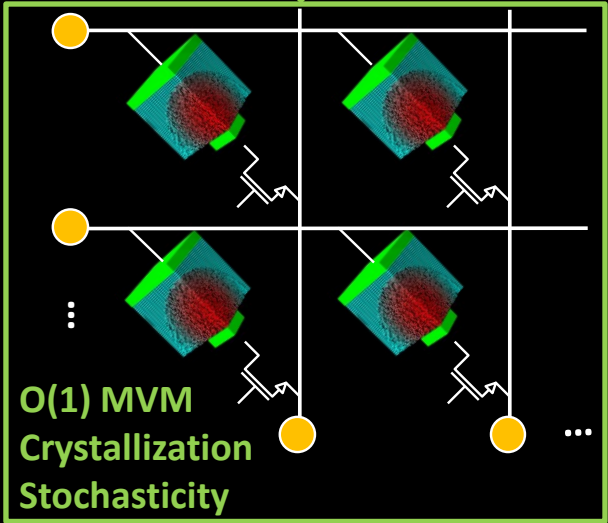
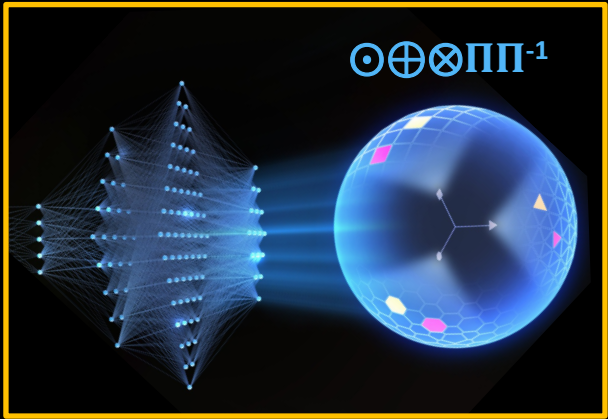


HW/SW realization is benefited from in-memory computing

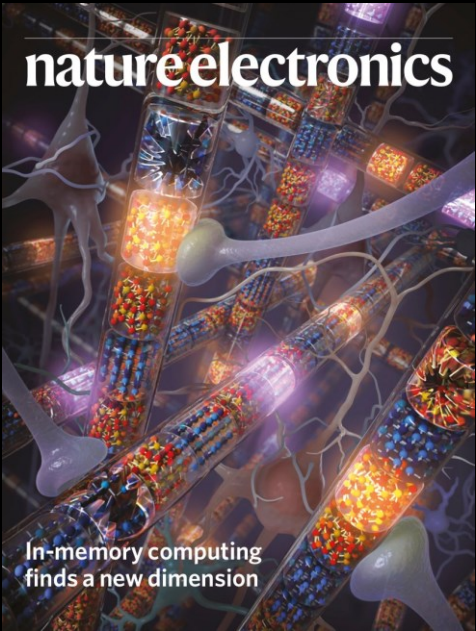


In-sensor few-shot learning

NVSA: From robust and efficient inference to learning and reasoning



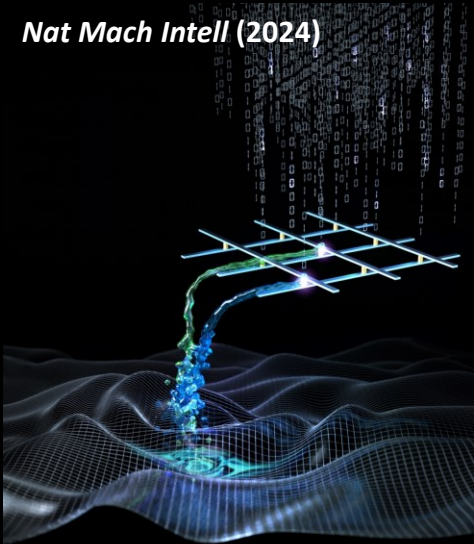
HW/SW realization is benefited from in-memory computing



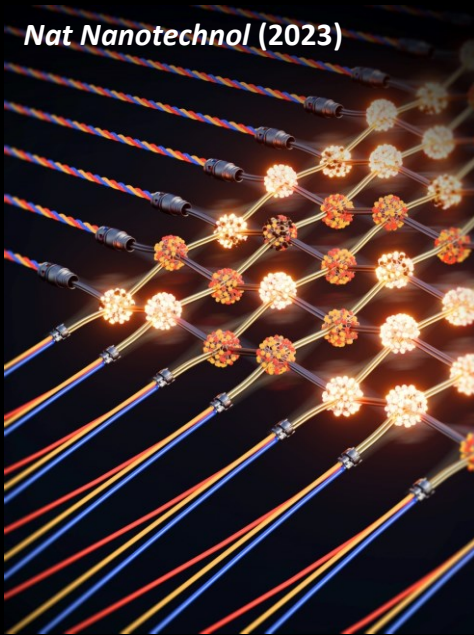
Robust VSA inference



In-sensor few-shot learning



In-memory kernel approximation

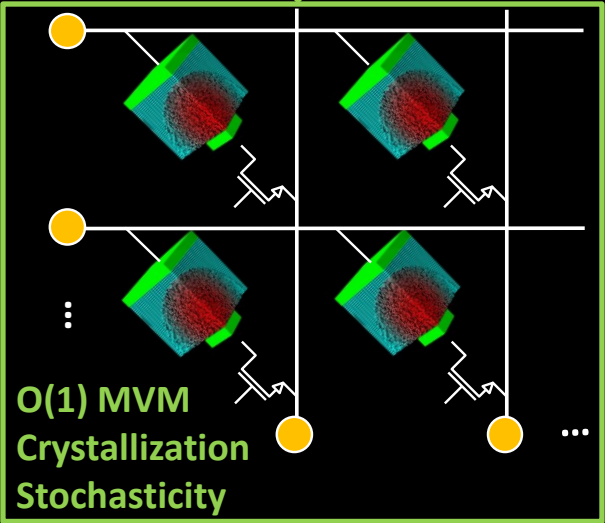
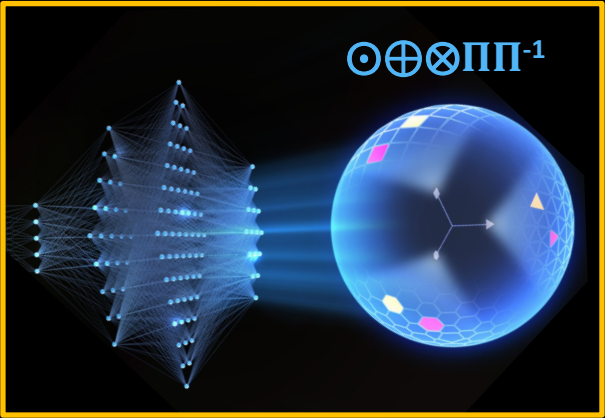


Disentangling representations

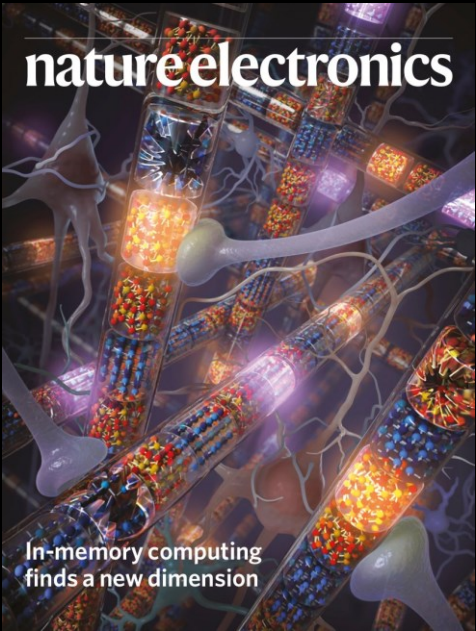


3D stack meets Mixture of Experts

NVSA: From robust and efficient inference to learning and reasoning



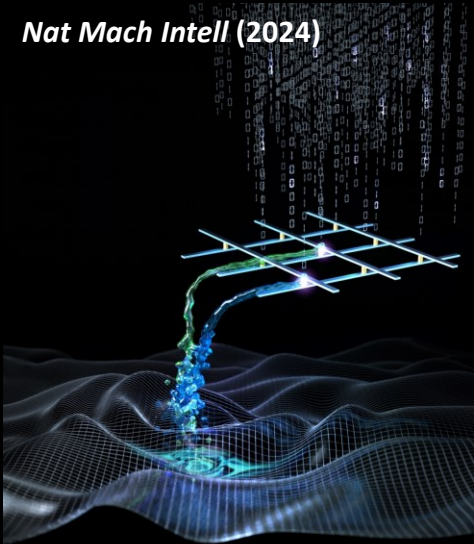
HW/SW realization is benefited from in-memory computing



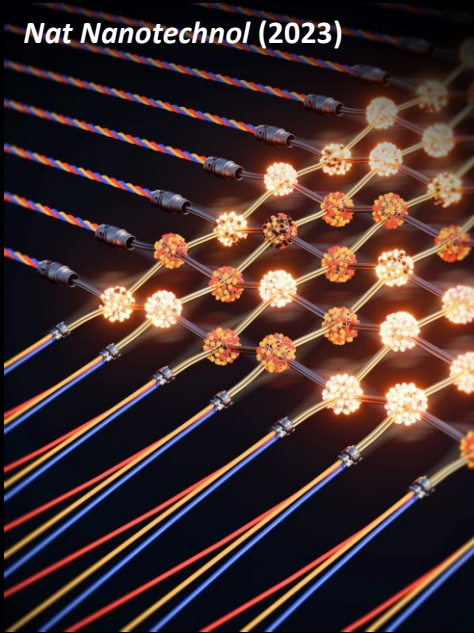
Robust VSA inference



In-sensor few-shot learning



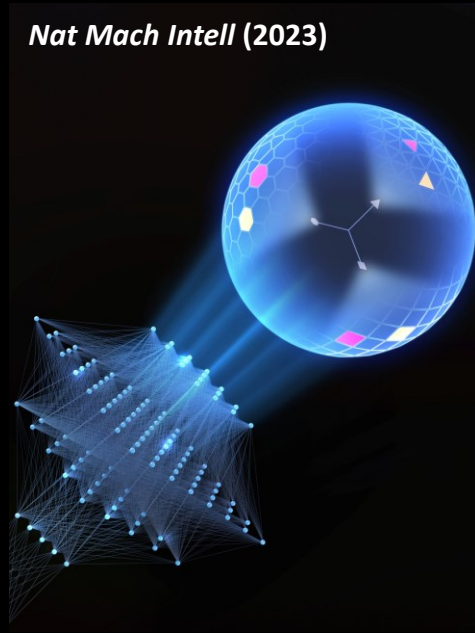
In-memory kernel approximation



Disentangling representations

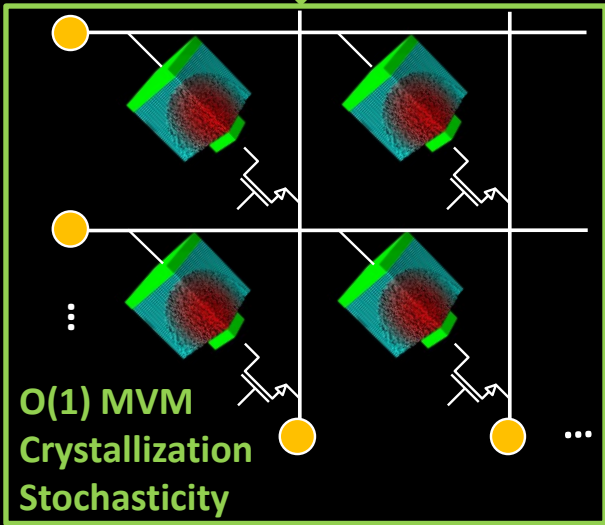
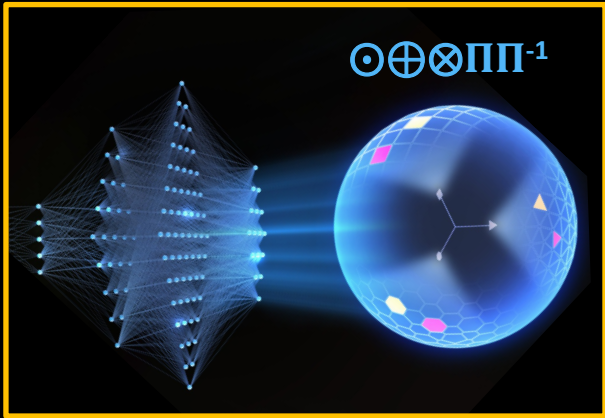


3D stack meets Mixture of Experts

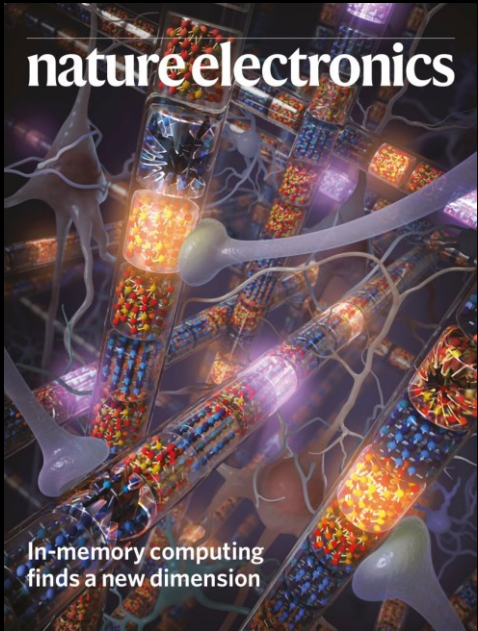


Scalable abstract reasoning

NVSA: From robust and efficient inference to learning and reasoning



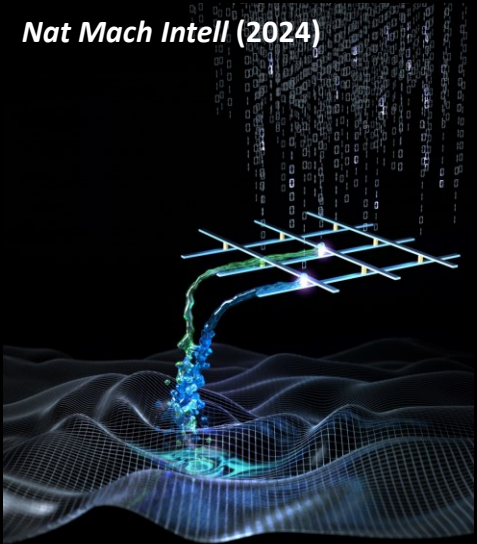
HW/SW realization is benefited from in-memory computing



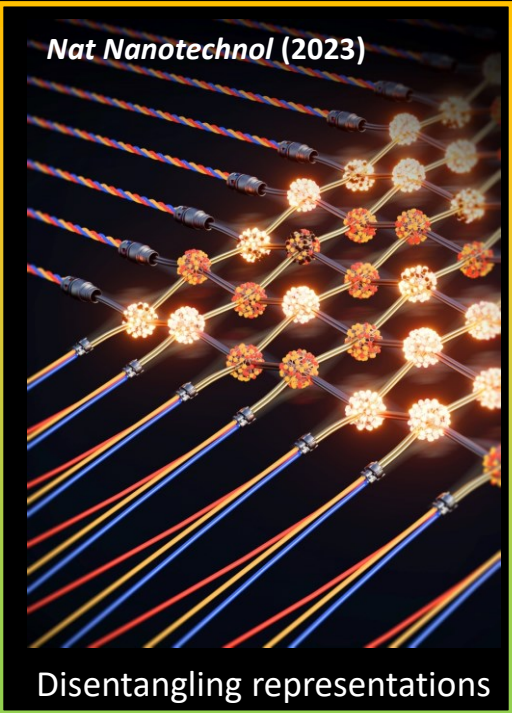
Robust VSA inference



In-sensor few-shot learning



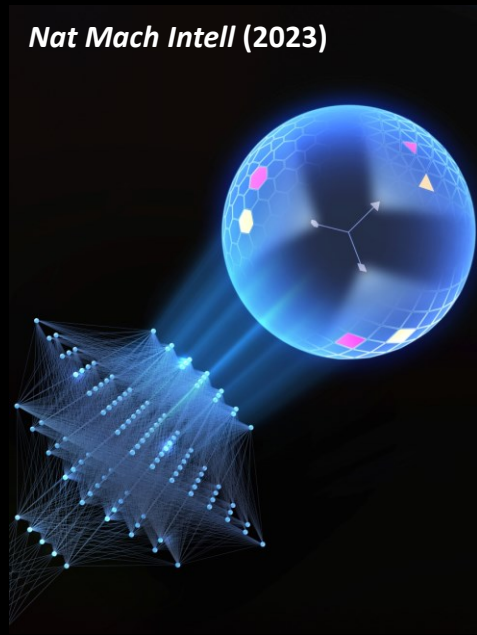
In-memory kernel approximation



Disentangling representations

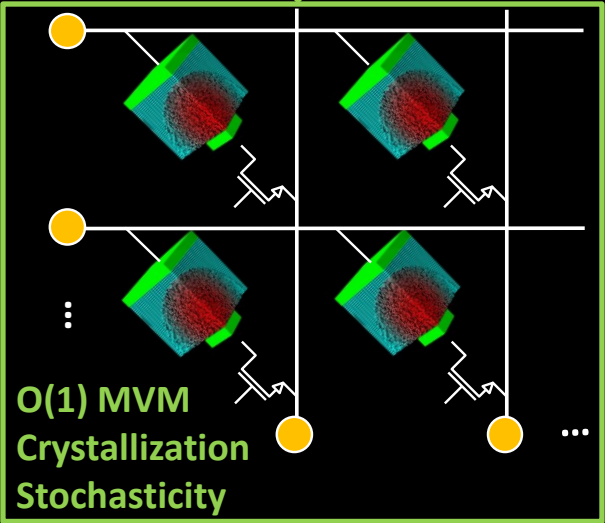
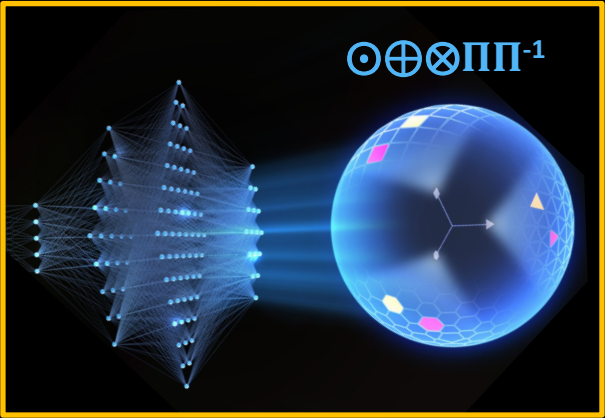


3D stack meets Mixture of Experts

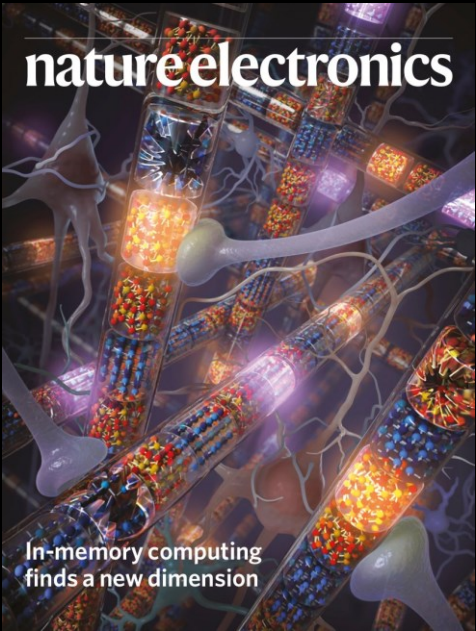


Scalable abstract reasoning

NVSA: From robust and efficient inference to learning and reasoning



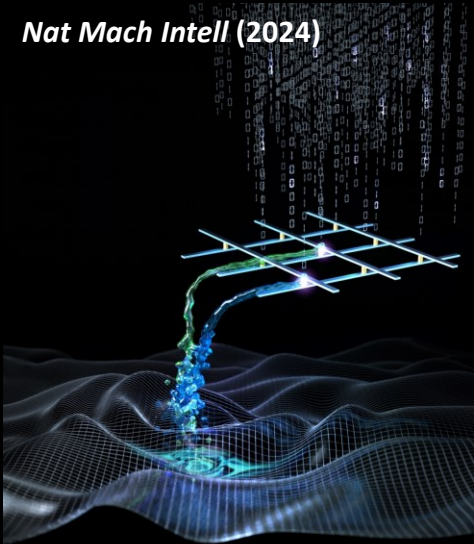
HW/SW realization is benefited from in-memory computing



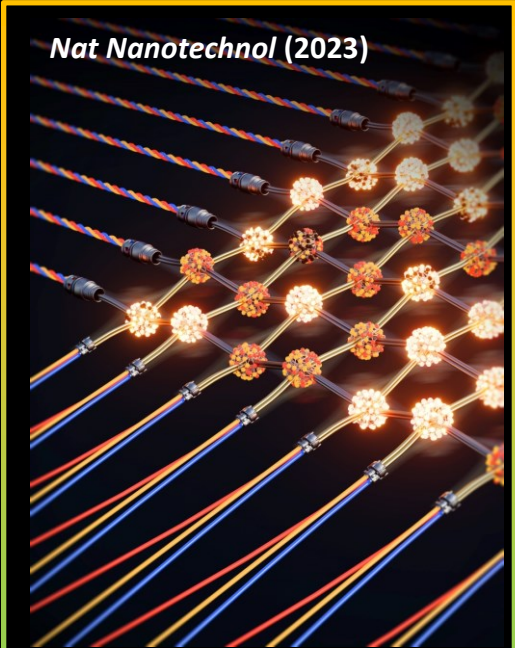
Robust VSA inference



In-sensor few-shot learning



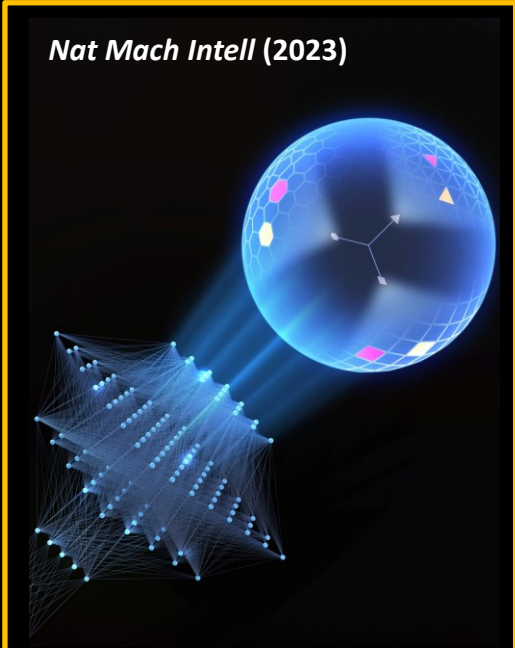
In-memory kernel approximation



Disentangling representations

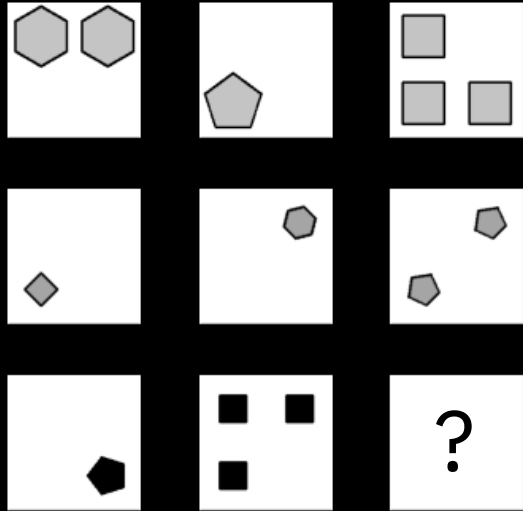


3D stack meets Mixture of Experts



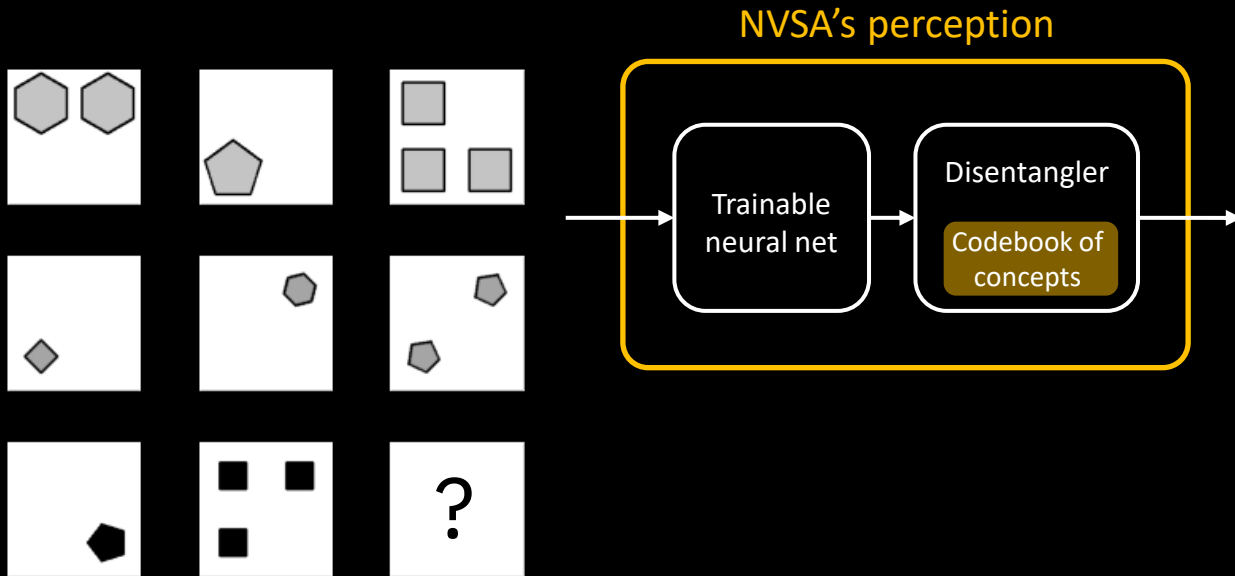
Scalable abstract reasoning

Visual abstract reasoning:



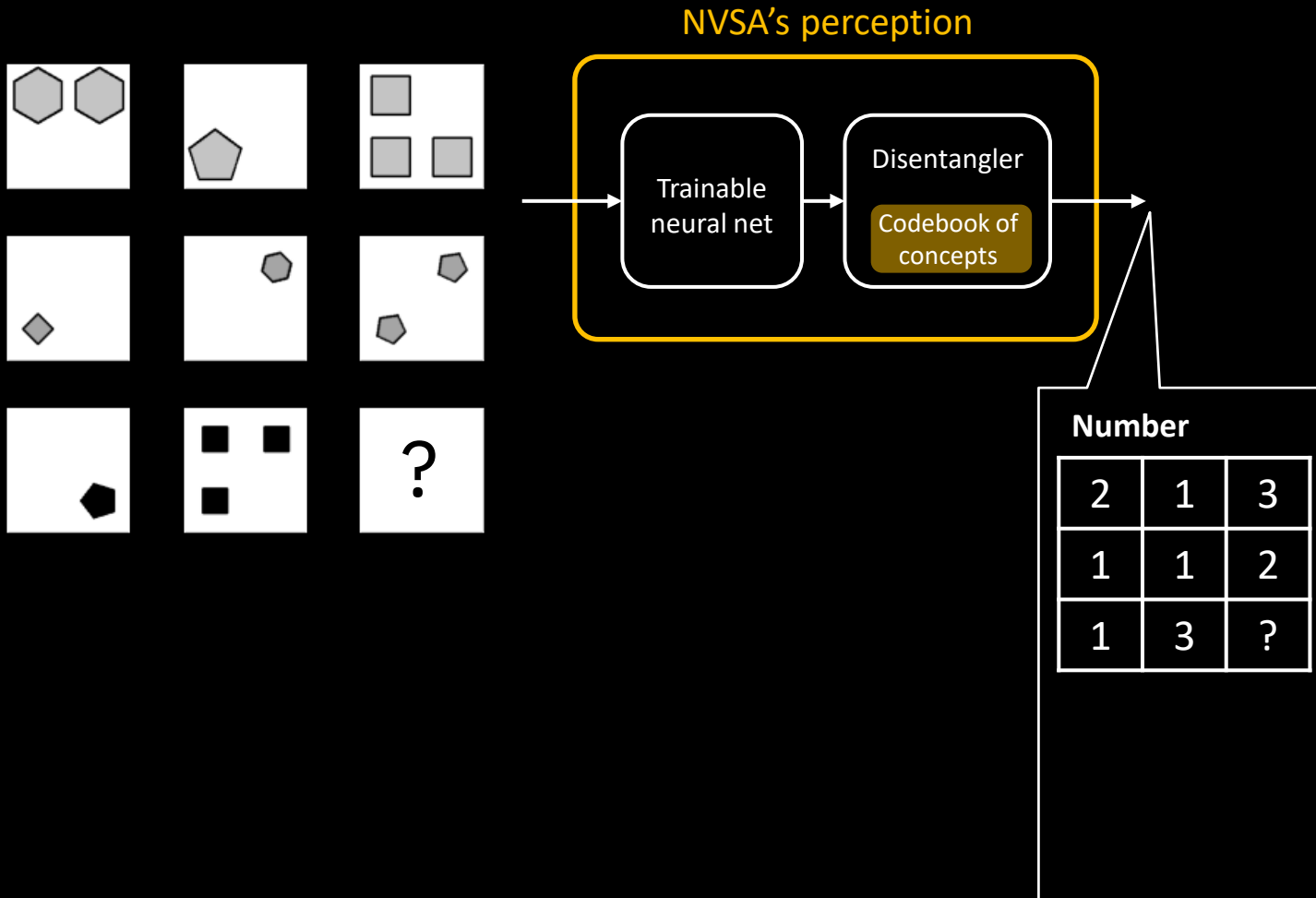
Visual abstract reasoning:

1) Disentangling perceptual representations



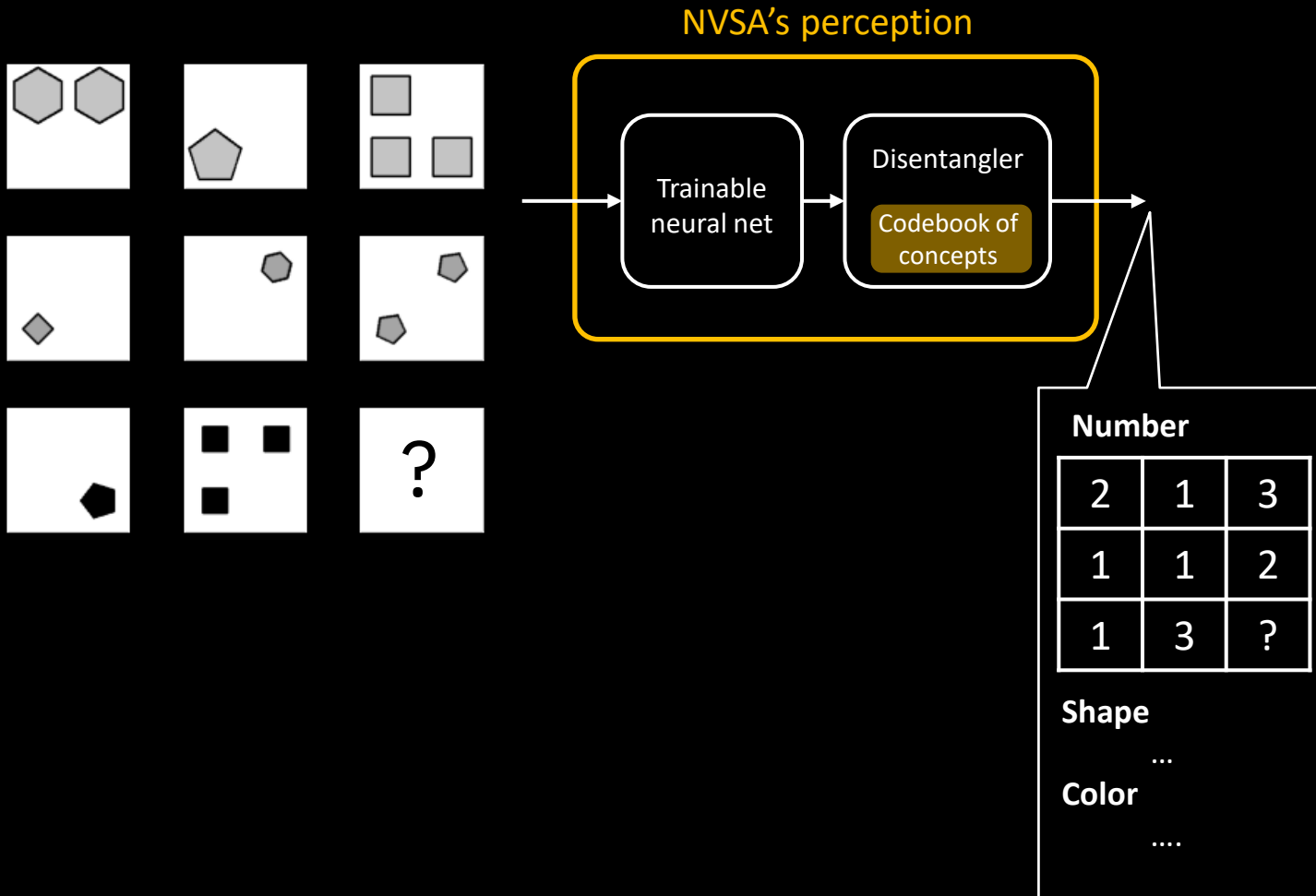
Visual abstract reasoning:

1) Disentangling perceptual representations



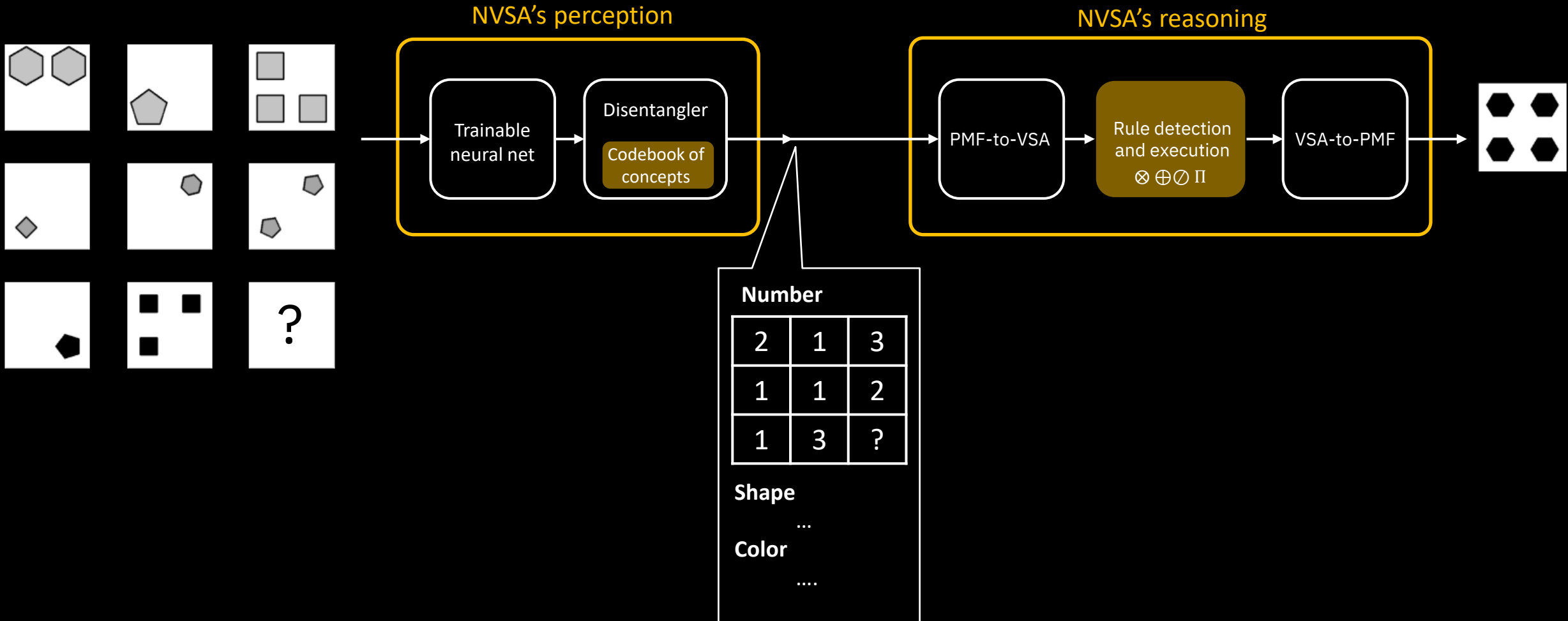
Visual abstract reasoning:

1) Disentangling perceptual representations



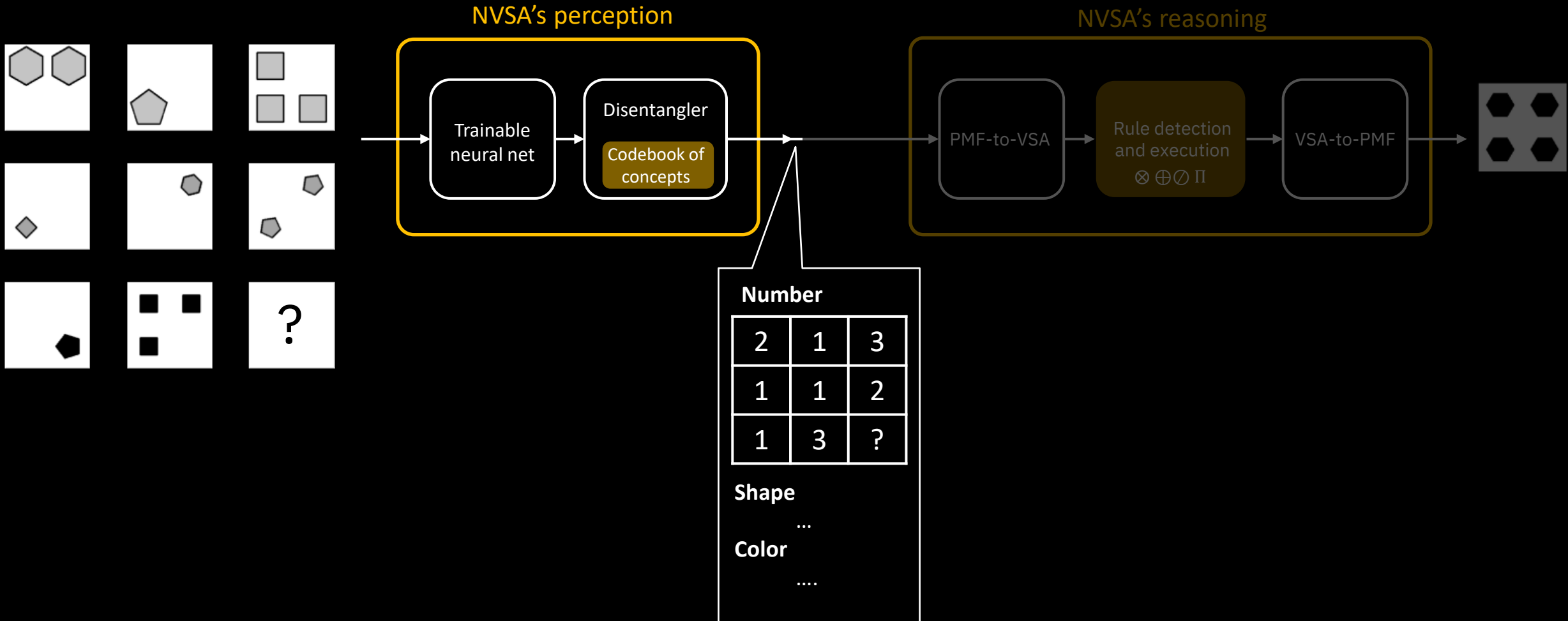
Visual abstract reasoning:

1) Disentangling perceptual representations

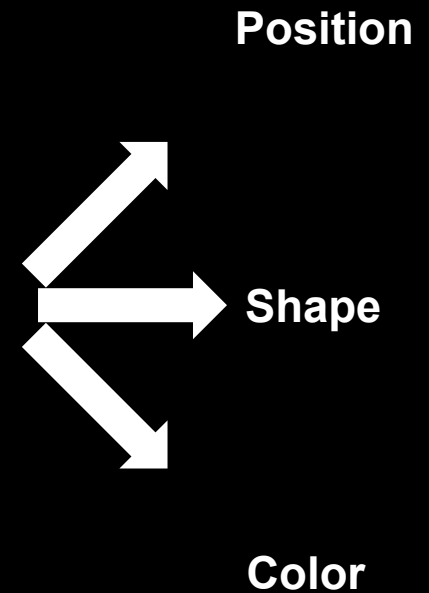
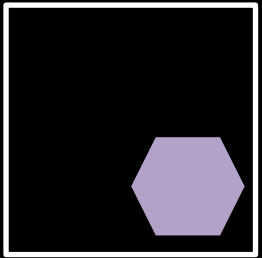


Visual abstract reasoning:

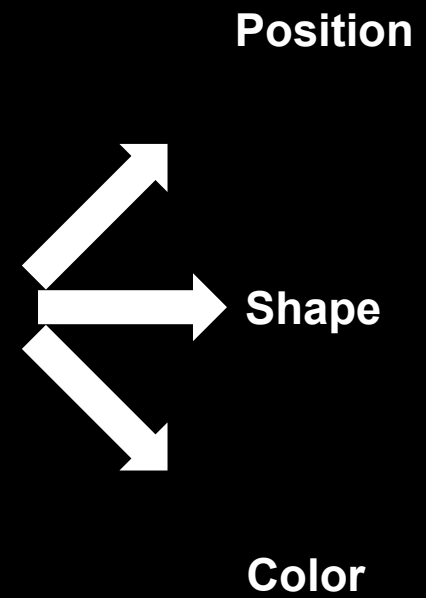
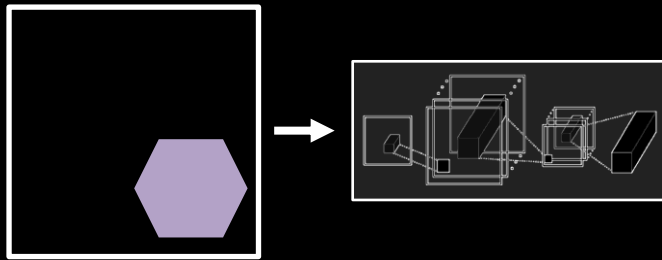
1) Disentangling perceptual representations



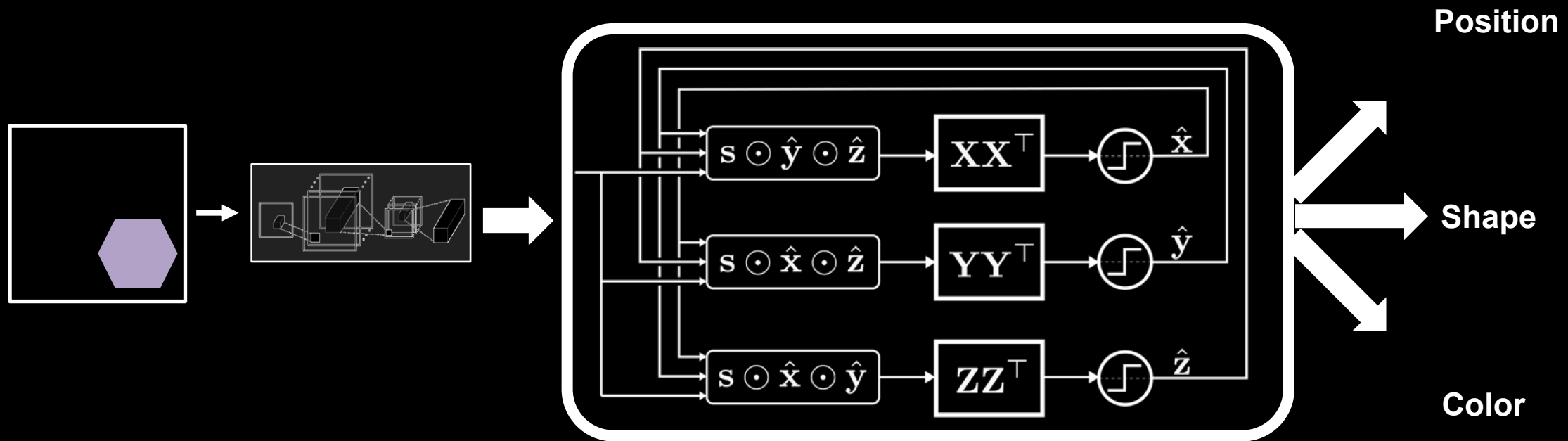
Efficiently disentangling visual representations: A dynamical system for searching in superposition



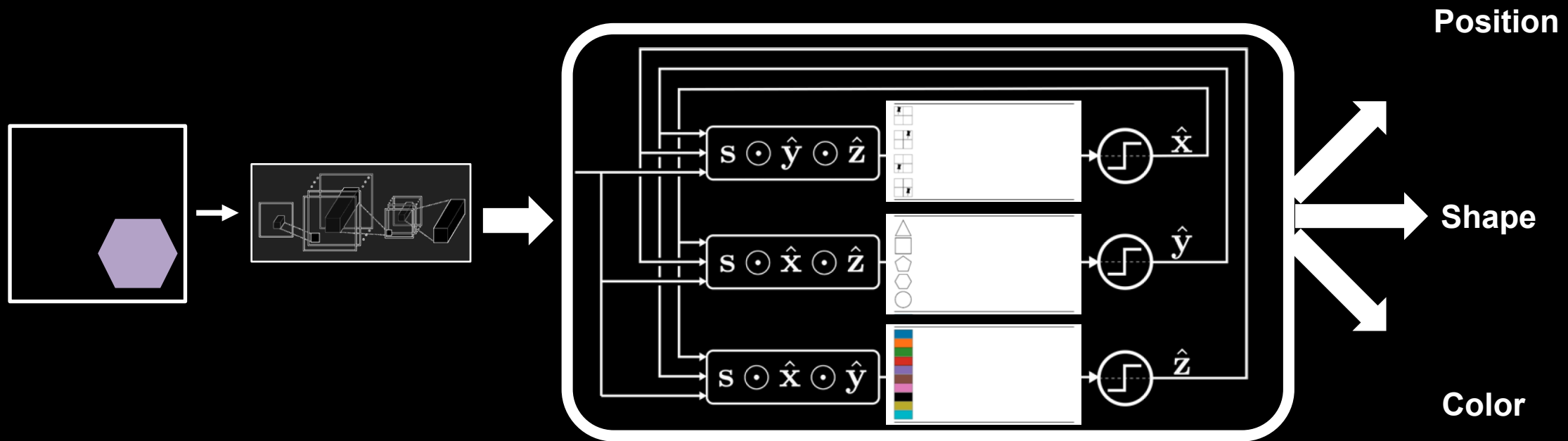
Efficiently disentangling visual representations: A dynamical system for searching in superposition



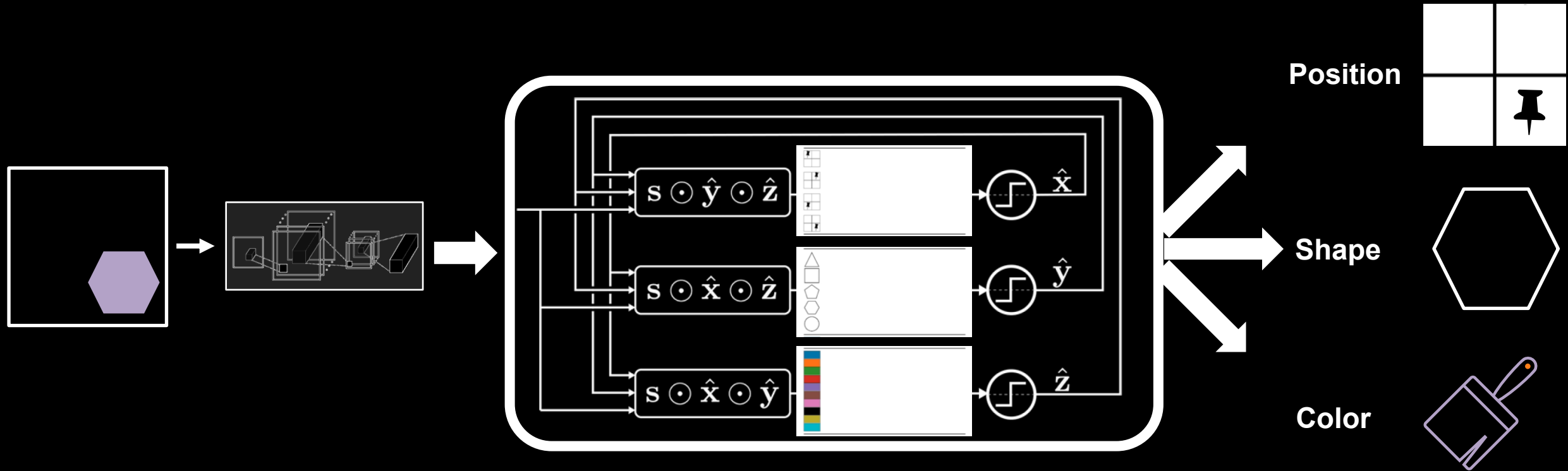
Efficiently disentangling visual representations: A dynamical system for searching in superposition



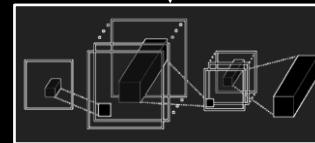
Efficiently disentangling visual representations: A dynamical system for searching in superposition



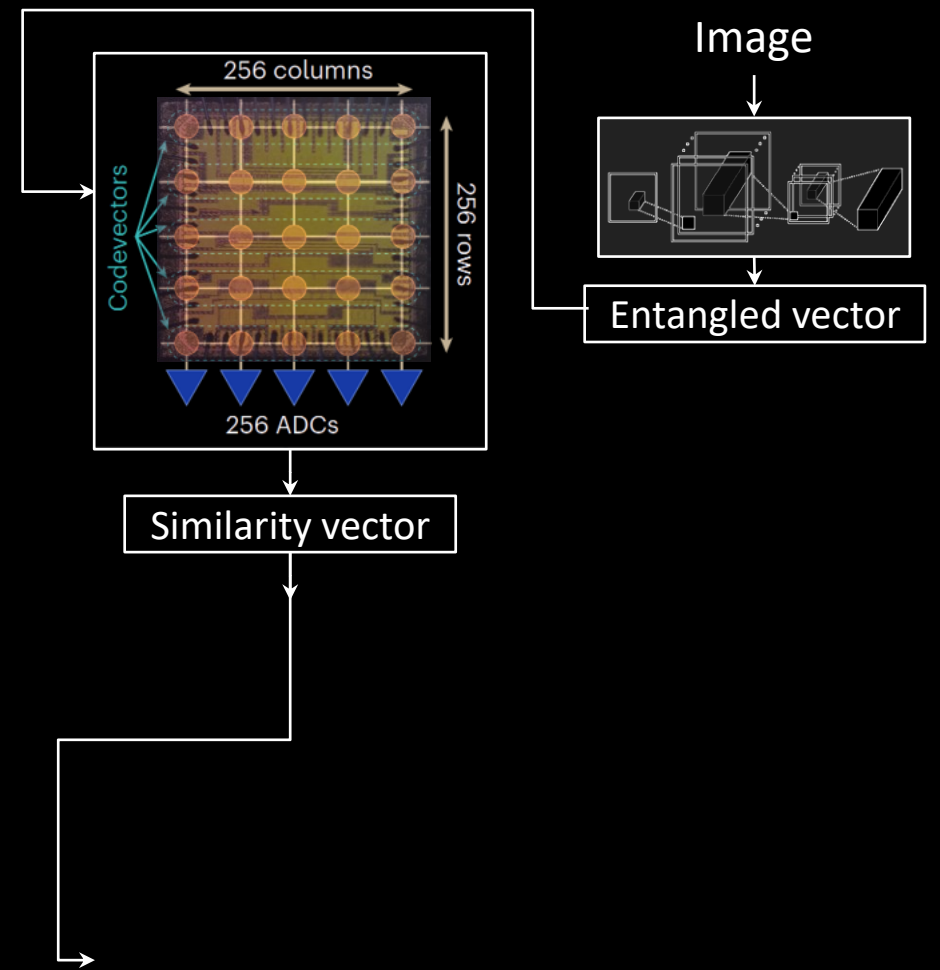
Efficiently disentangling visual representations: A dynamical system for searching in superposition

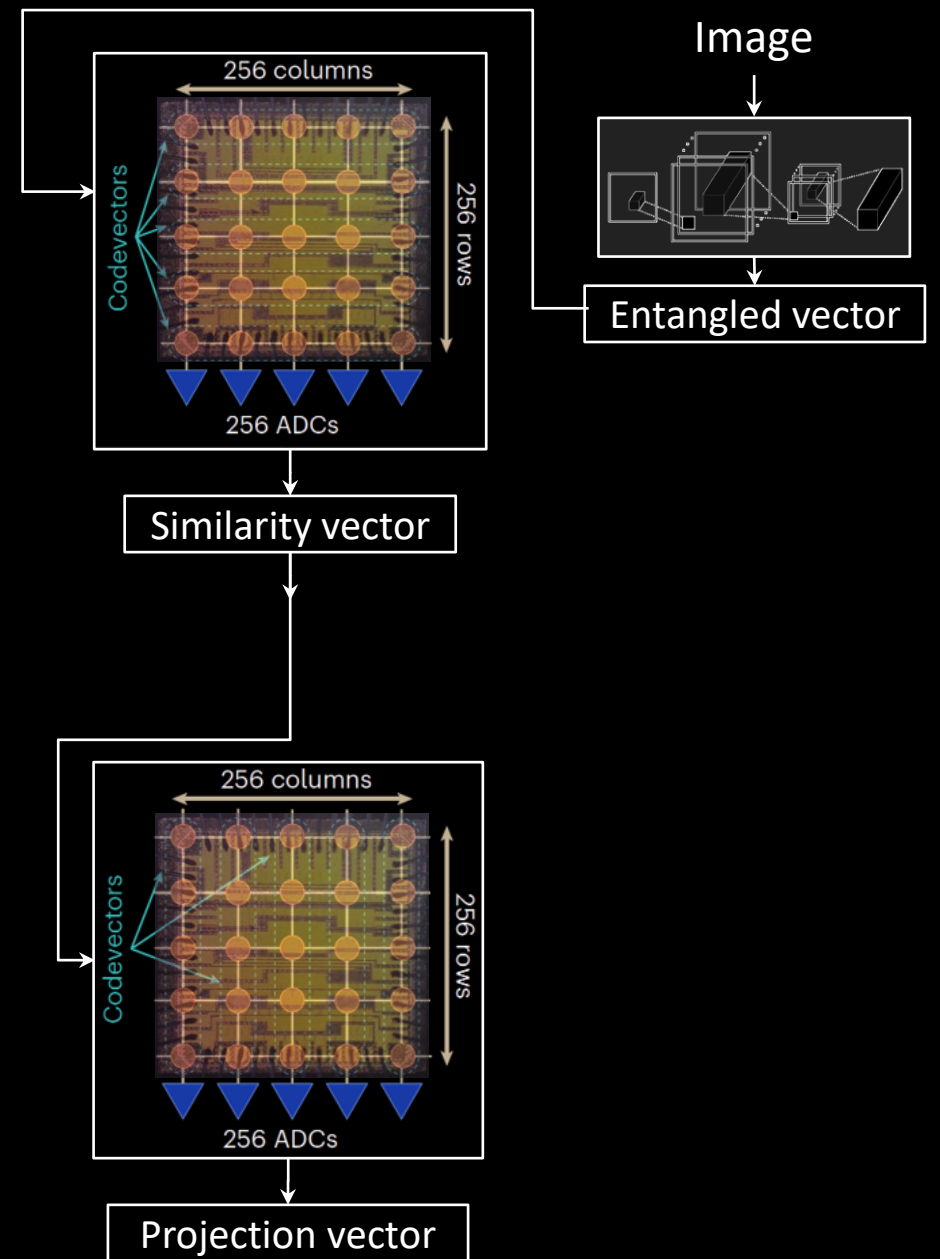


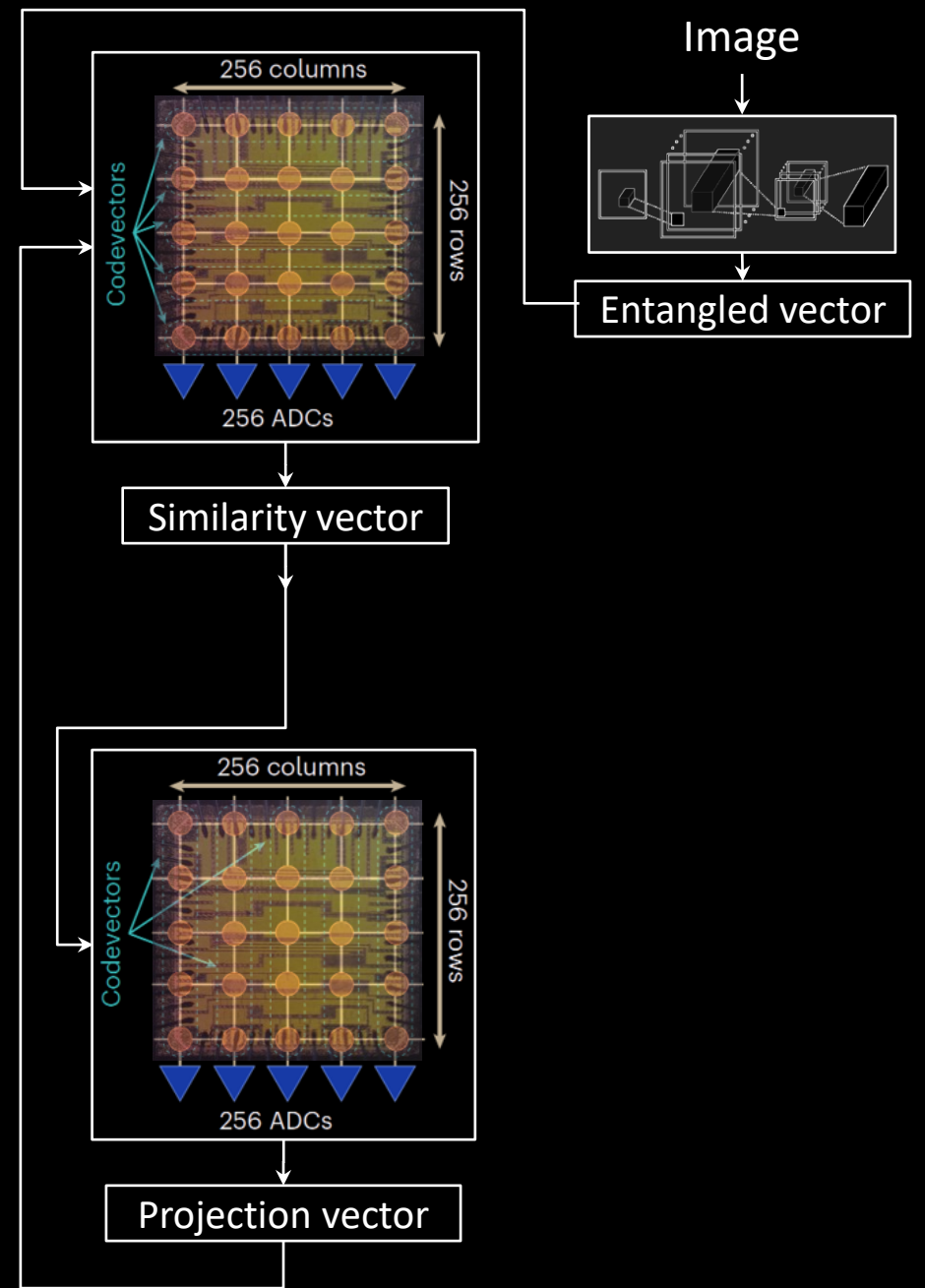
Image



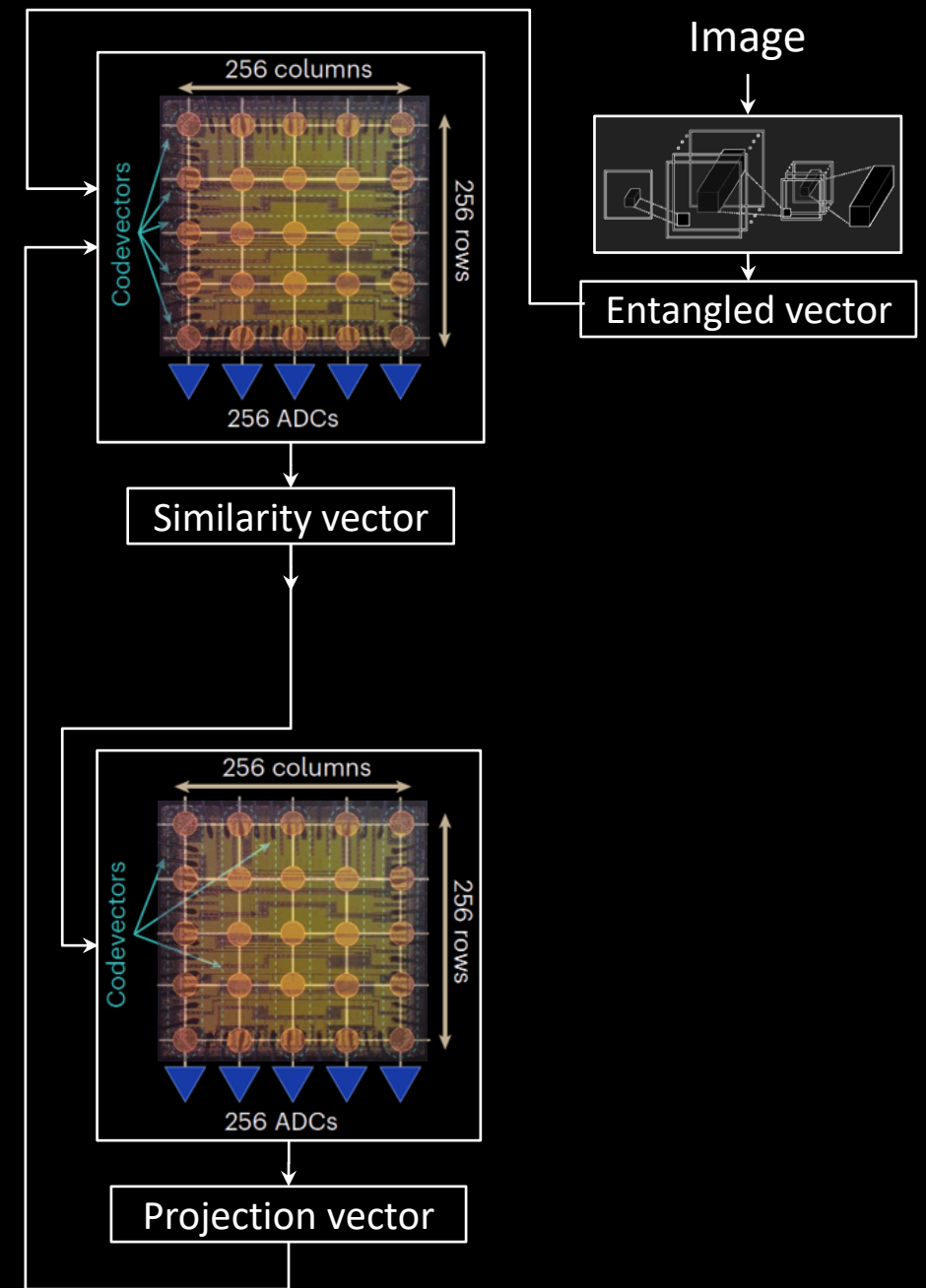
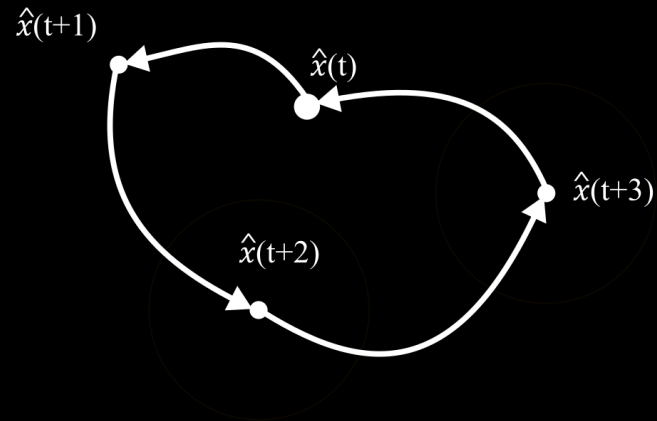
Entangled vector



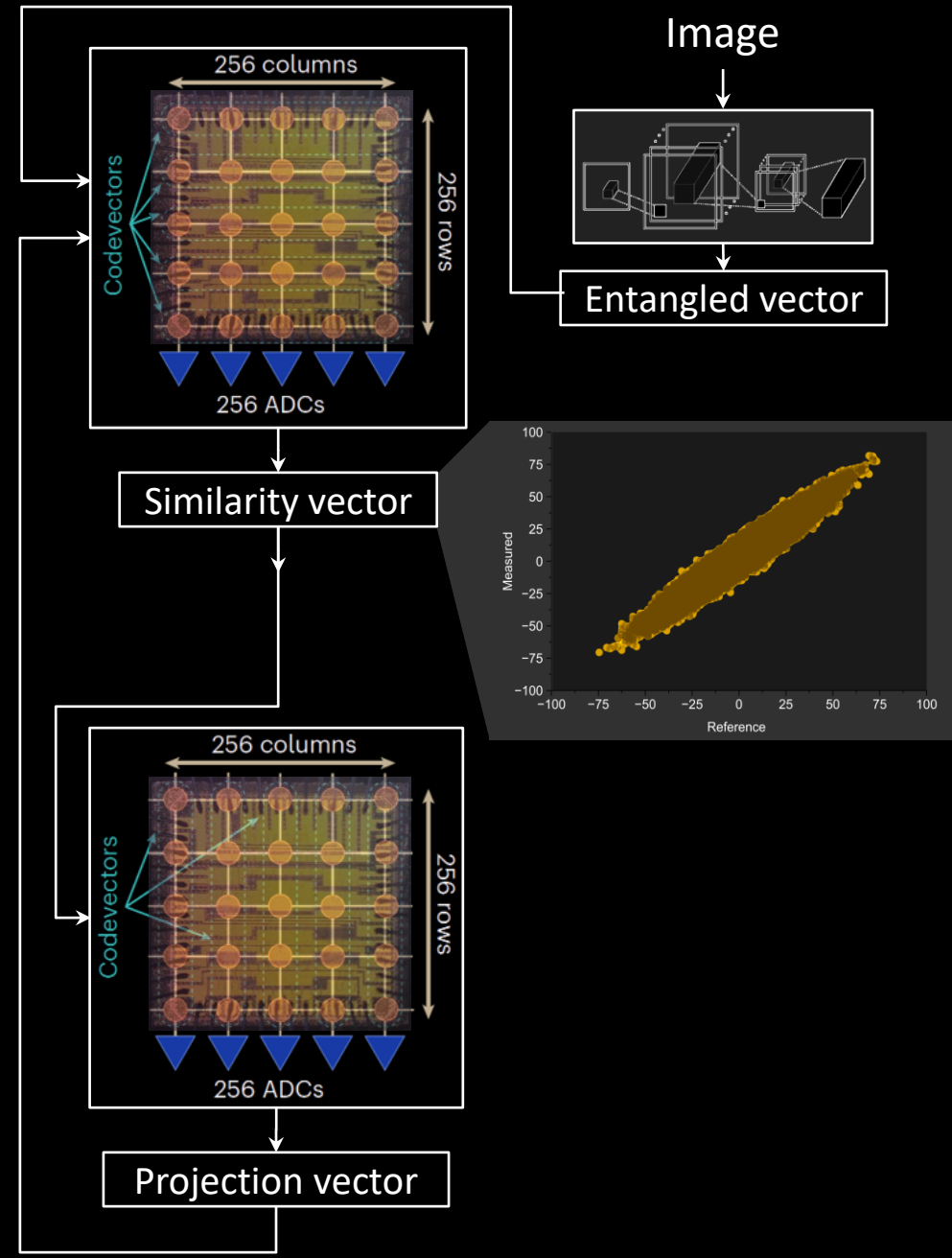
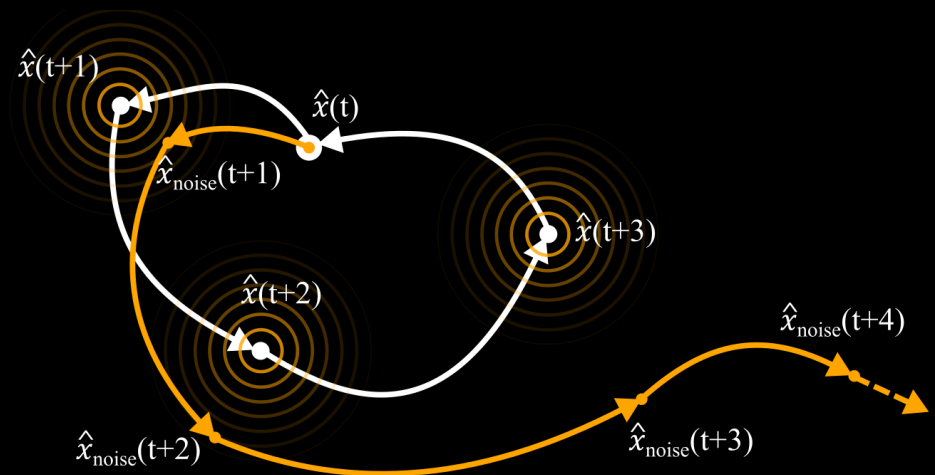




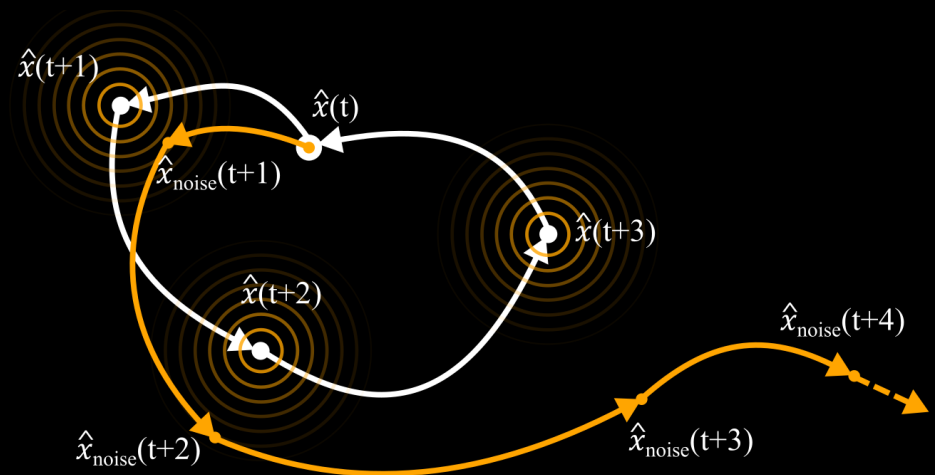
1) Faces limit cycles:



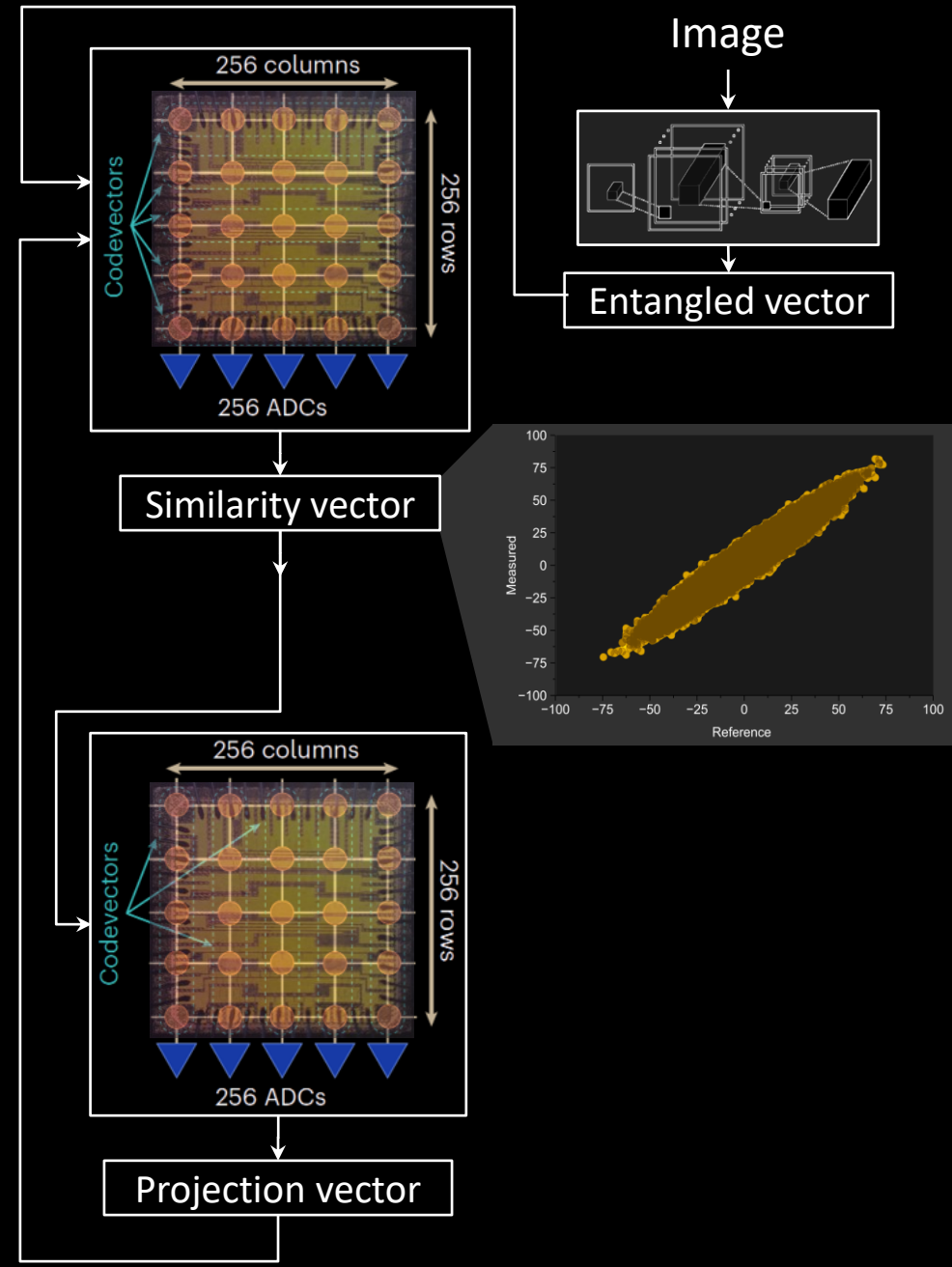
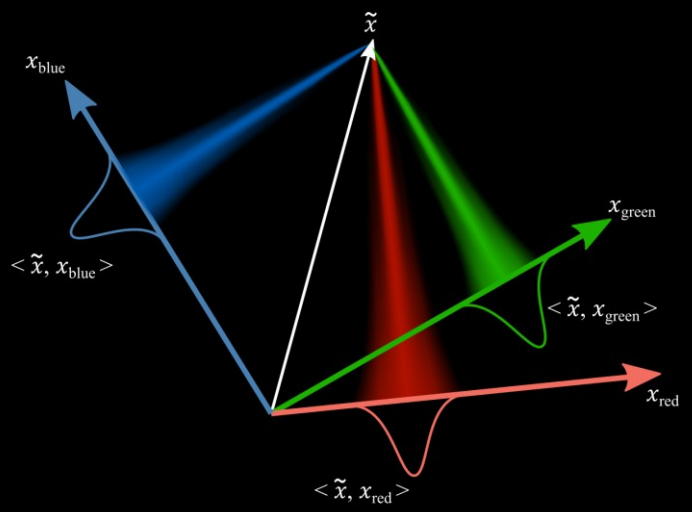
1) Faces limit cycles: removed by stochasticity



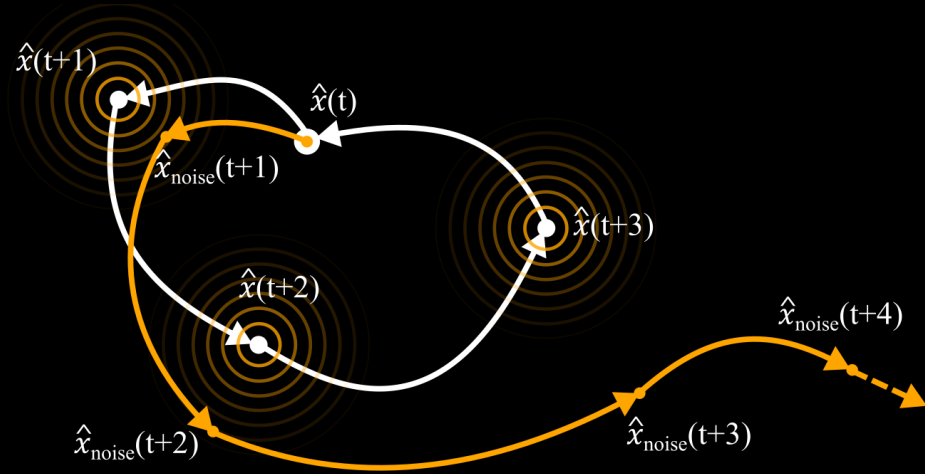
1) Faces limit cycles: removed by stochasticity



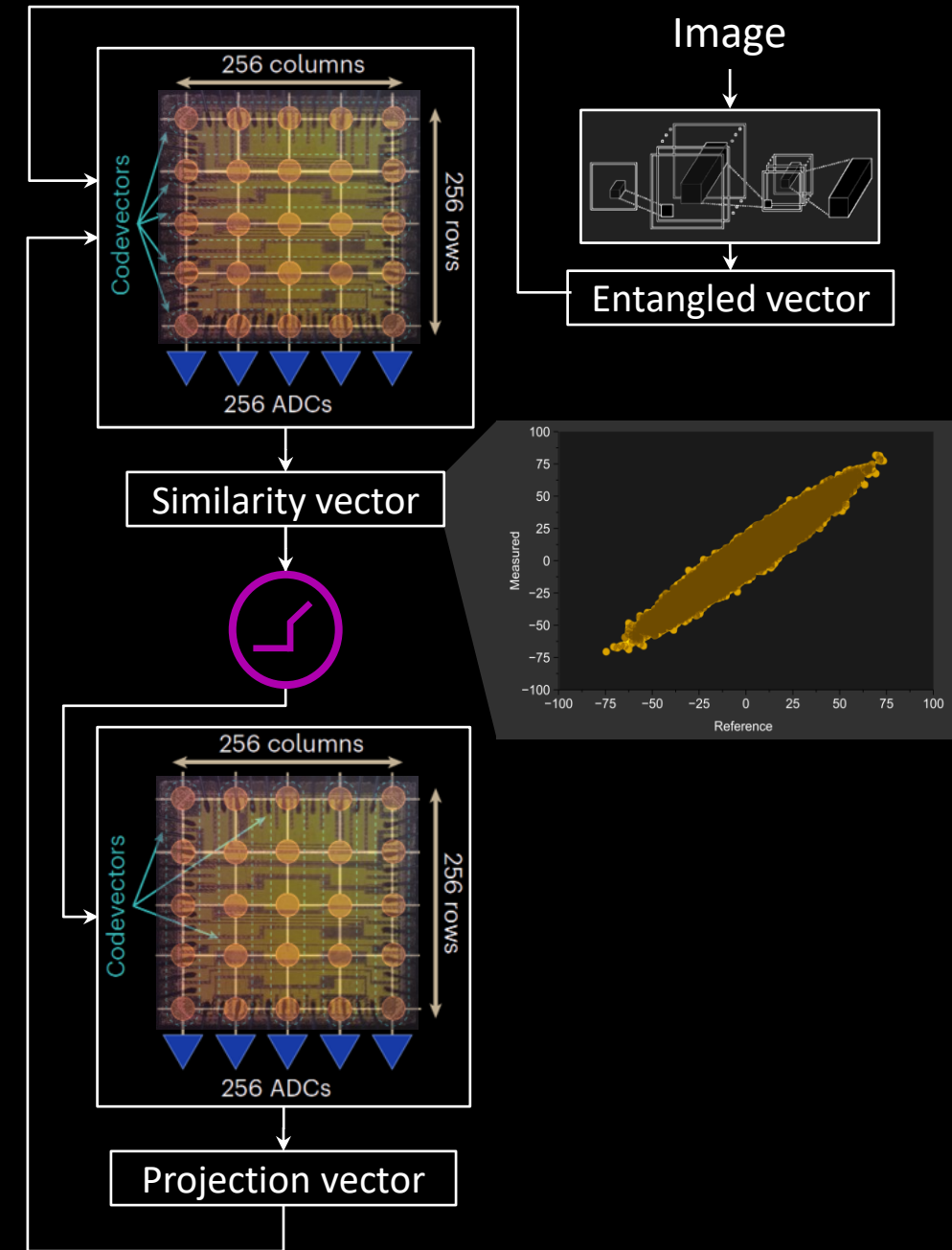
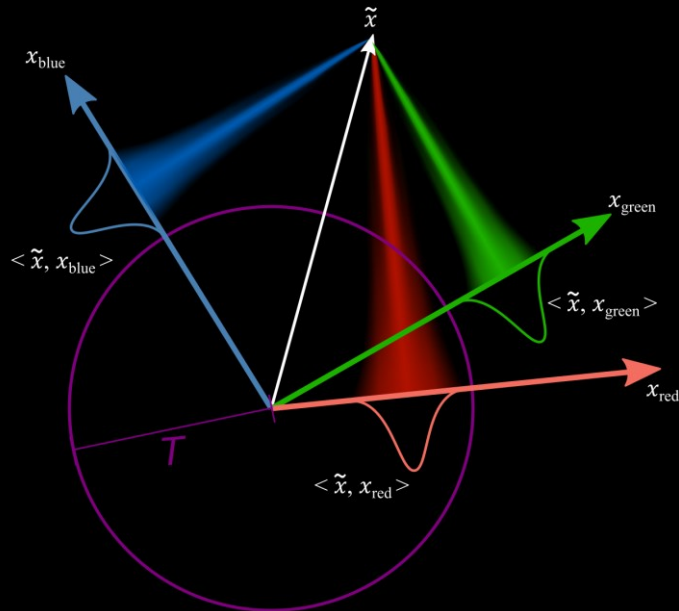
2) Limited capacity:

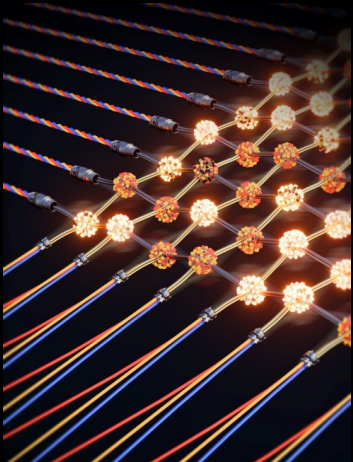


1) Faces limit cycles: removed by stochasticity



2) Limited capacity: boosted by sparse activation

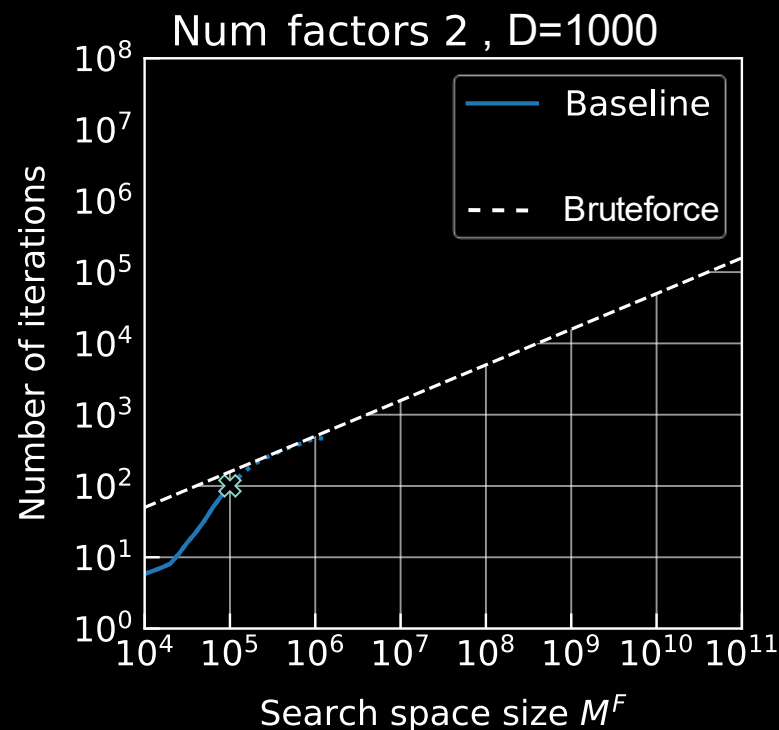




Stochastic in-memory factorizer (IMF) improves capacity by at least **five orders of magnitude**

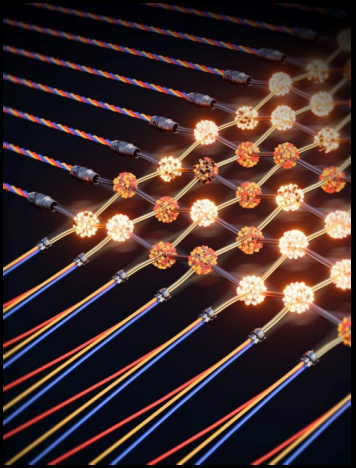
10× faster convergence

12× lower energy



[Langenegger et al., In-memory factorization of holographic perceptual representations, *Nature Nanotechnology*, 2023]

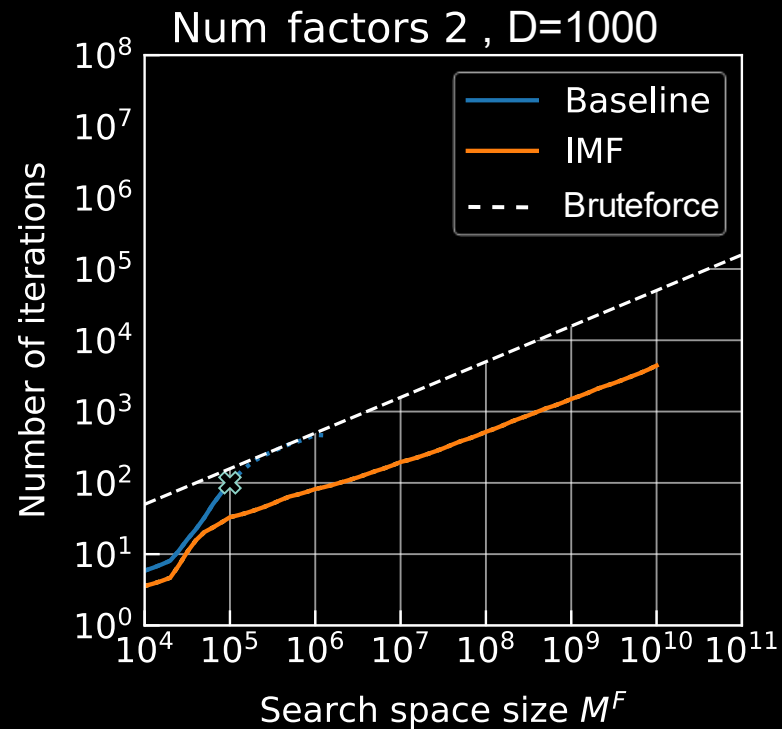
Highlighted in Nat. Nanotechnol.'s News & Views



Stochastic in-memory factorizer (IMF) improves capacity by at least **five orders of magnitude**

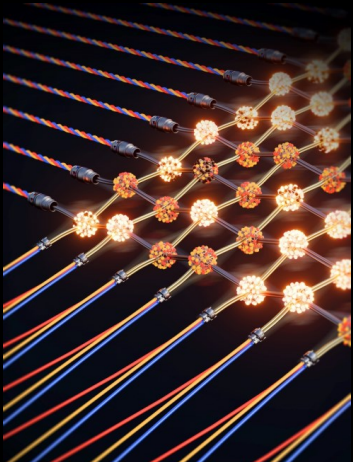
10× faster convergence

12× lower energy



[Langenegger et al., In-memory factorization of holographic perceptual representations, *Nature Nanotechnology*, 2023]

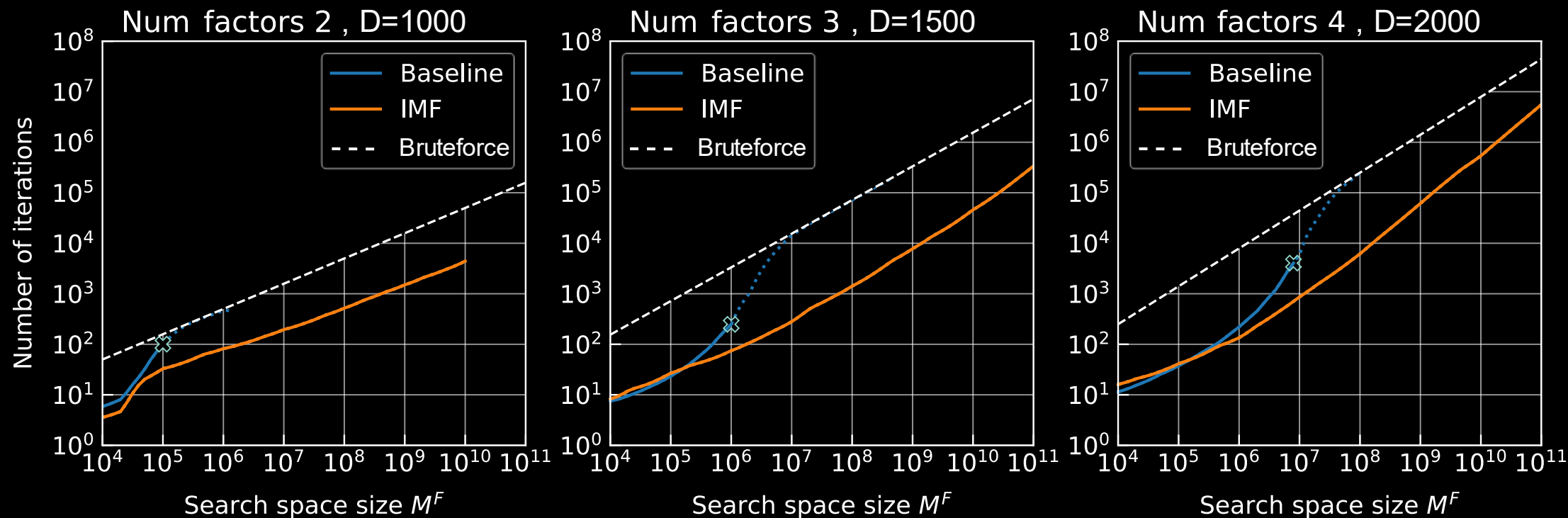
Highlighted in Nat. Nanotechnol.'s News & Views



Stochastic in-memory factorizer (IMF) improves capacity by at least **five orders of magnitude**

10× faster convergence

12× lower energy

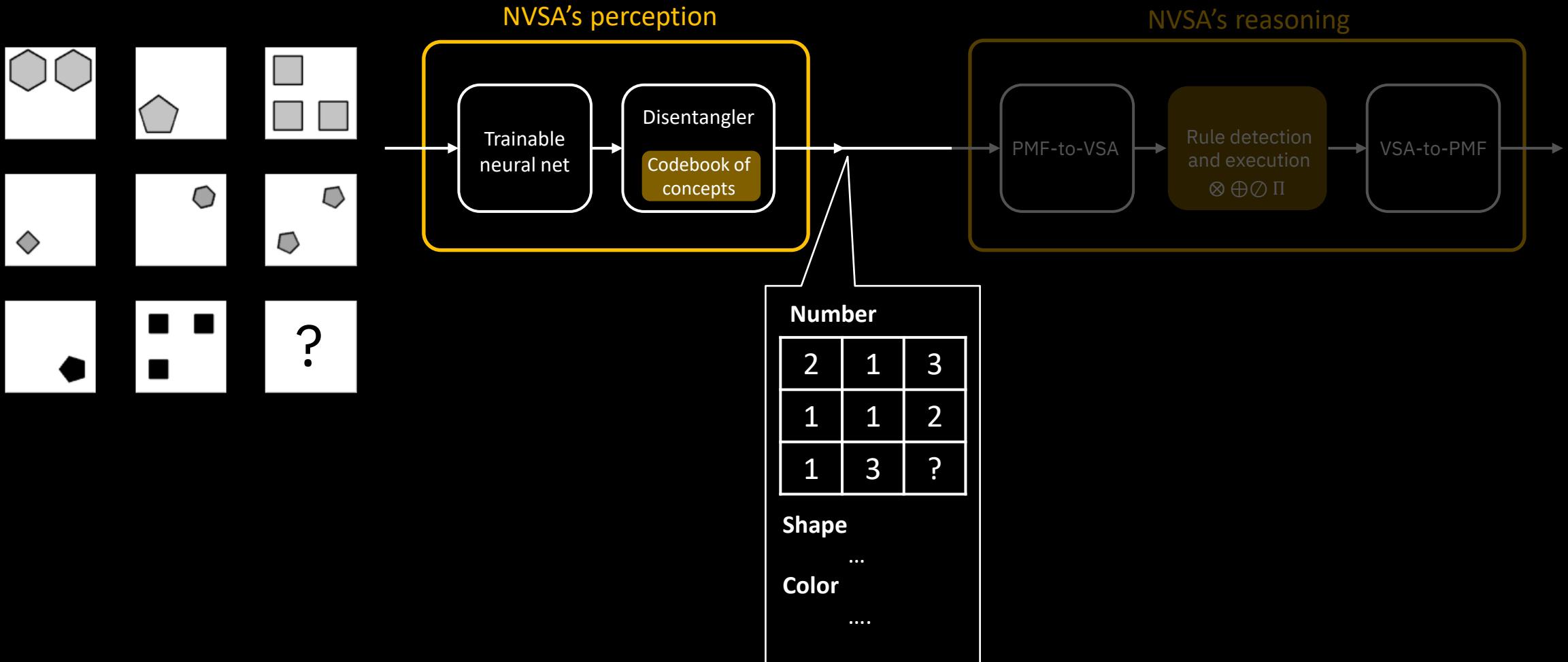


[Langenegger et al., In-memory factorization of holographic perceptual representations, *Nature Nanotechnology*, 2023]

Highlighted in Nat. Nanotechnol.'s News & Views

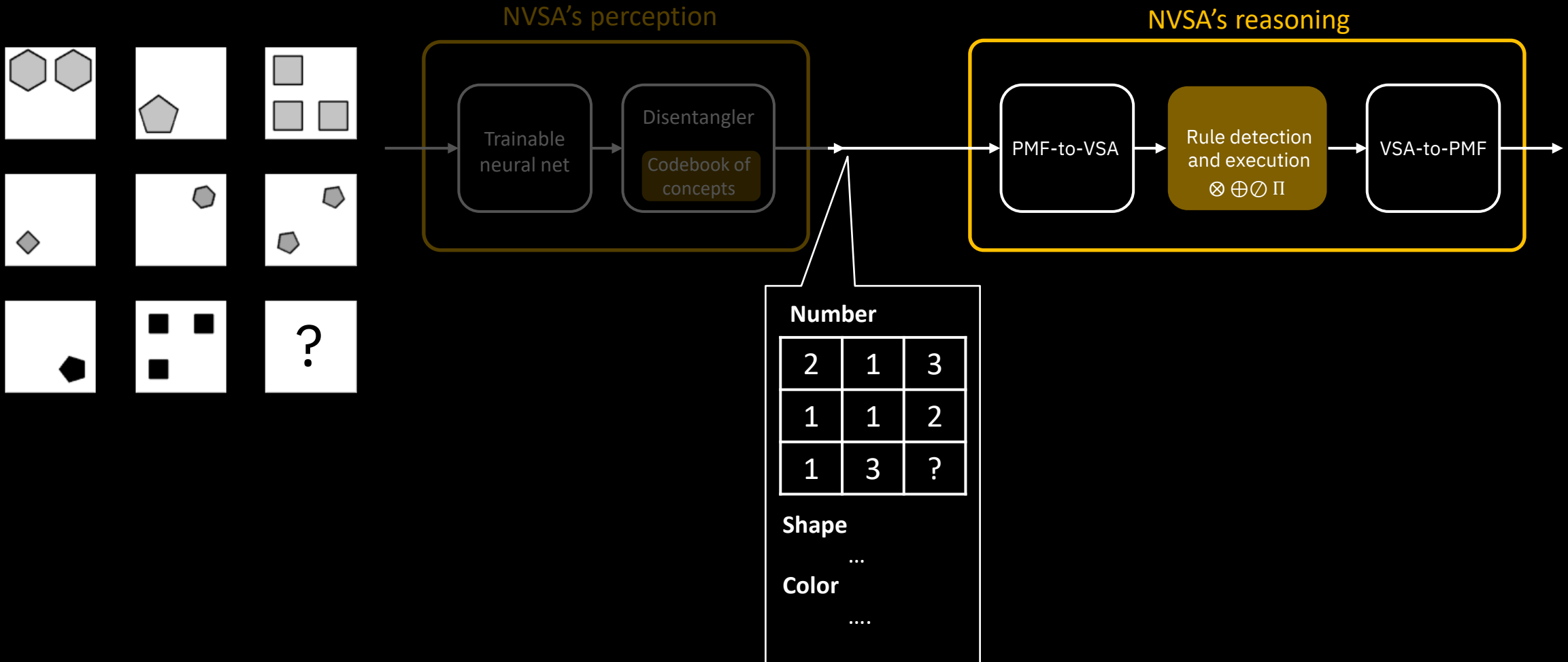
Visual abstract reasoning:

2) Scalable reasoning with disentangled representations



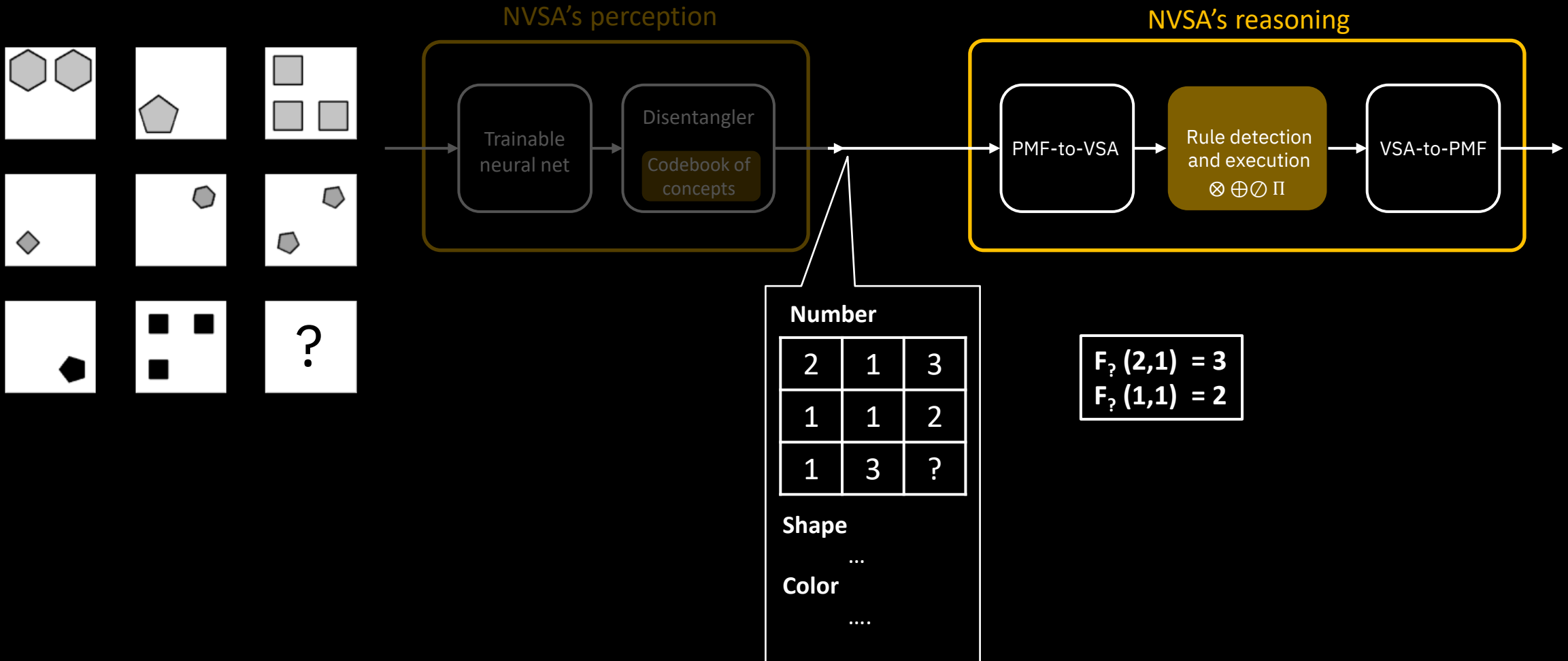
Visual abstract reasoning:

2) Scalable reasoning with disentangled representations



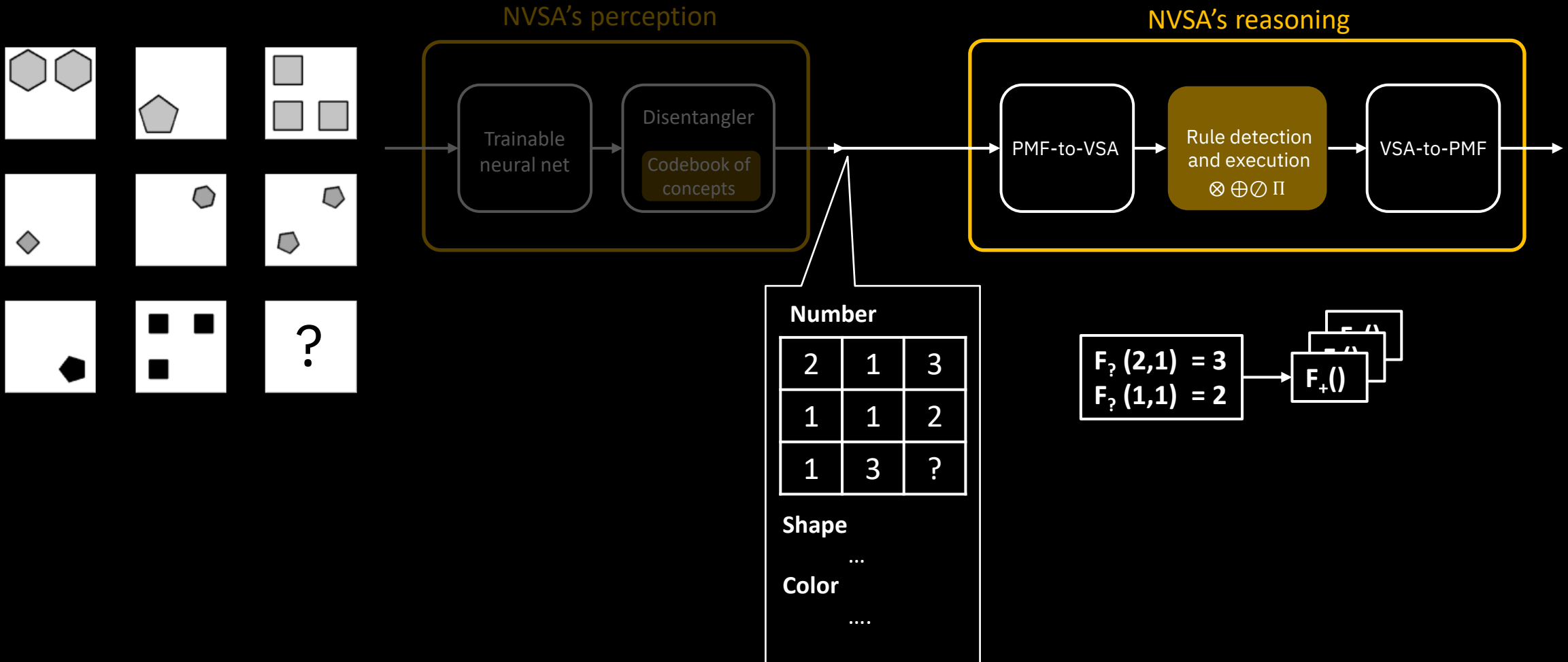
Visual abstract reasoning:

2) Scalable reasoning with disentangled representations

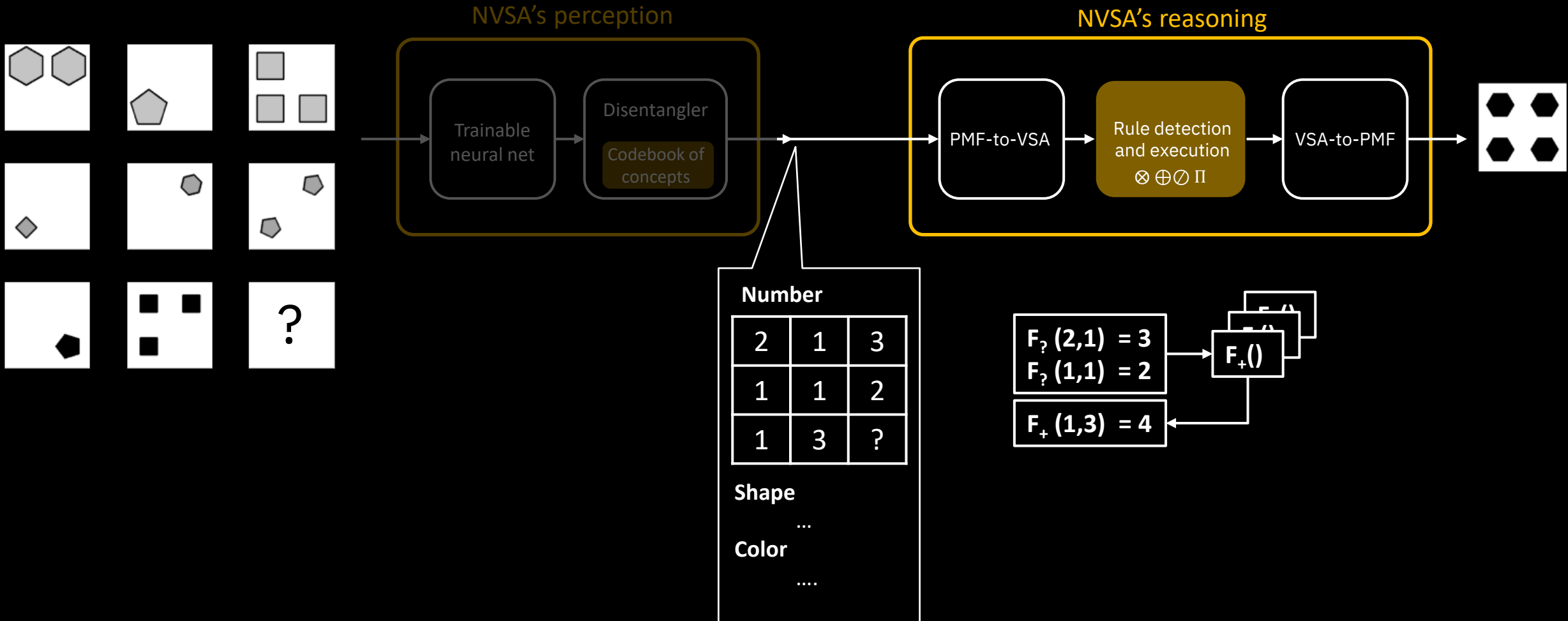


Visual abstract reasoning:

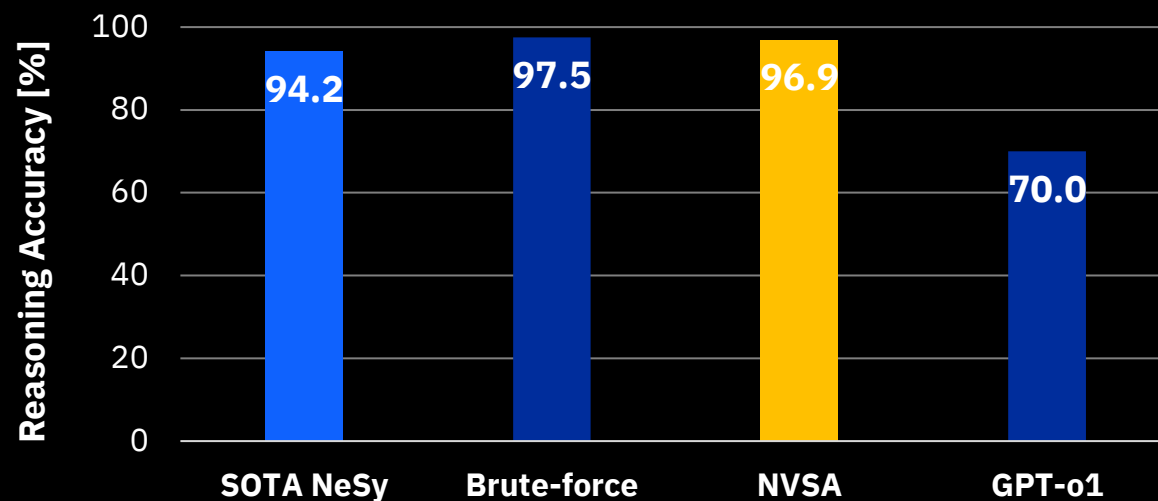
2) Scalable reasoning with disentangled representations



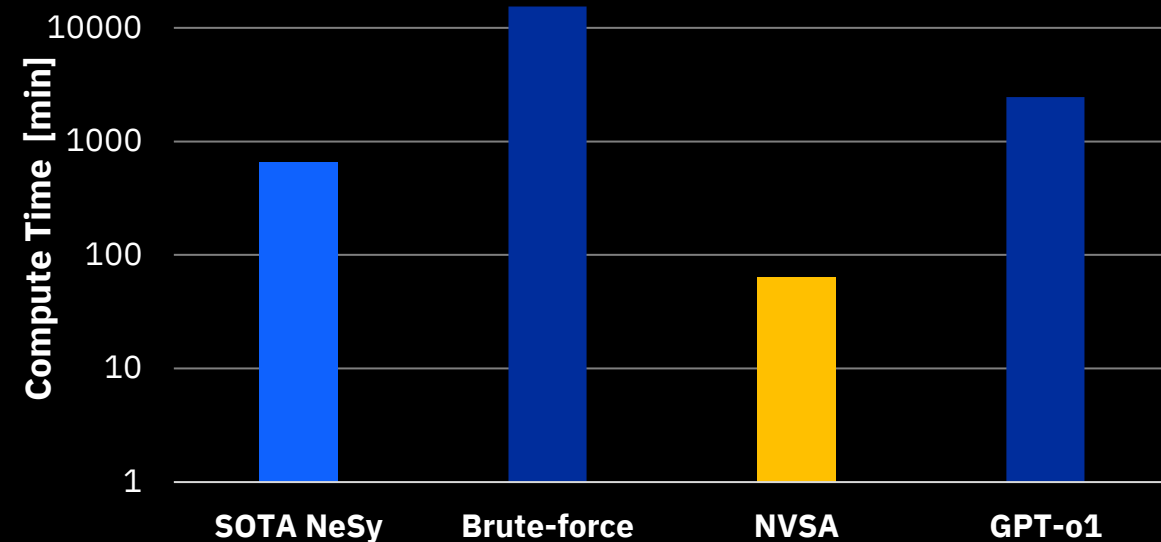
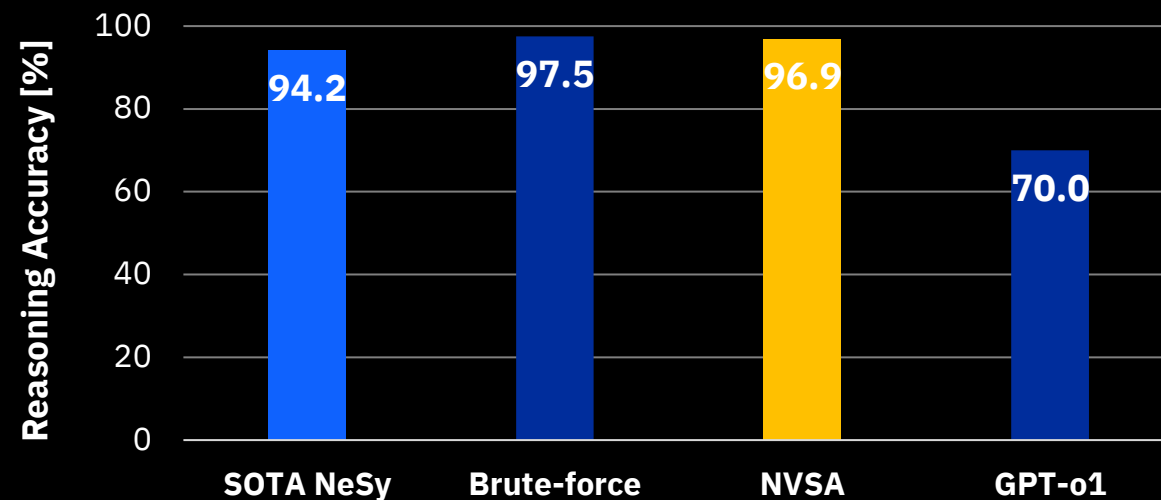
2) Scalable reasoning with disentangled representations



NVSA's reasoning is **240x** faster
enabling real-time inference on CPUs

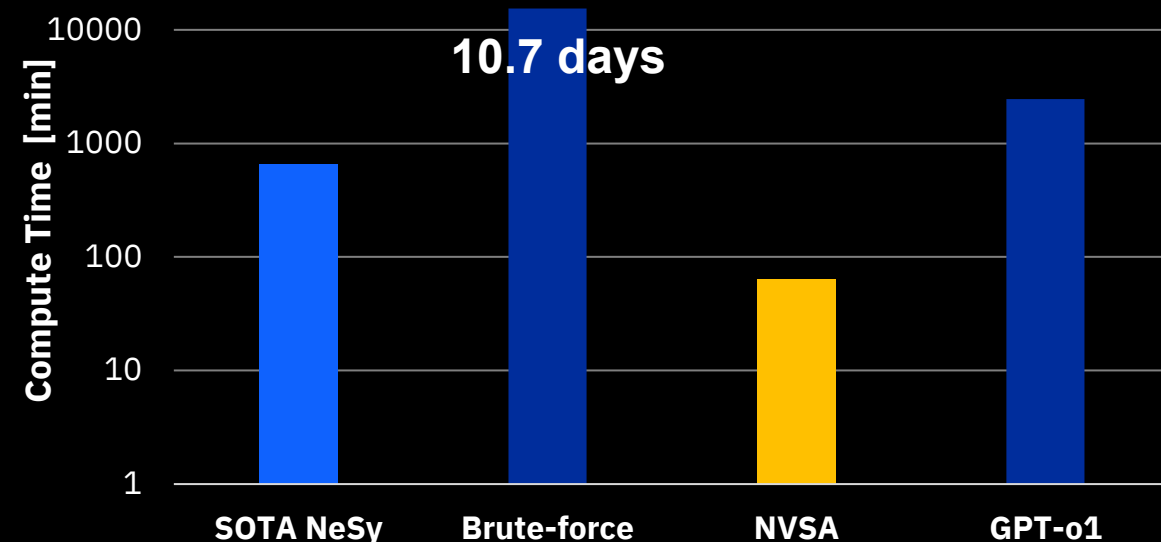
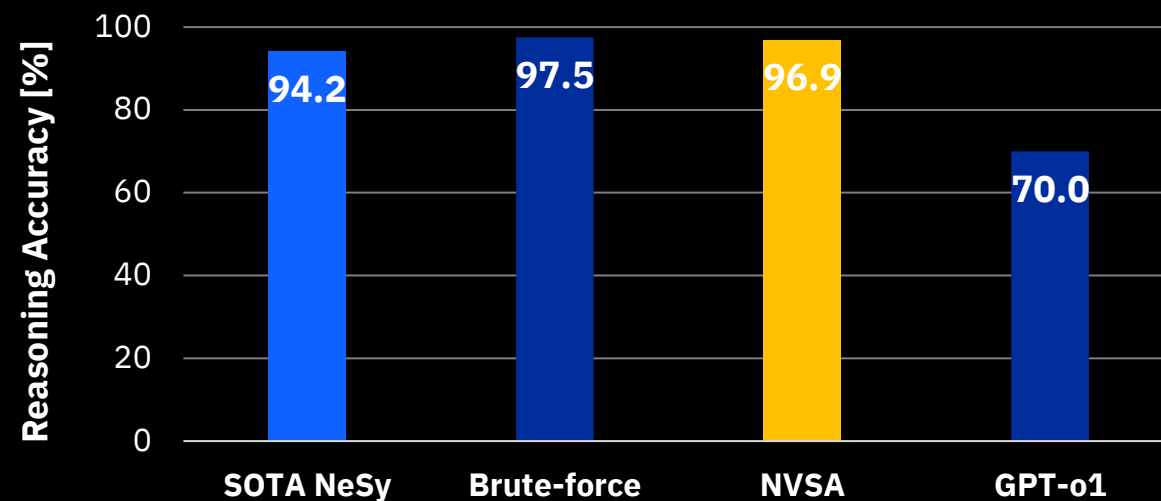


NVSA's reasoning is **240x** faster
enabling real-time inference on CPUs



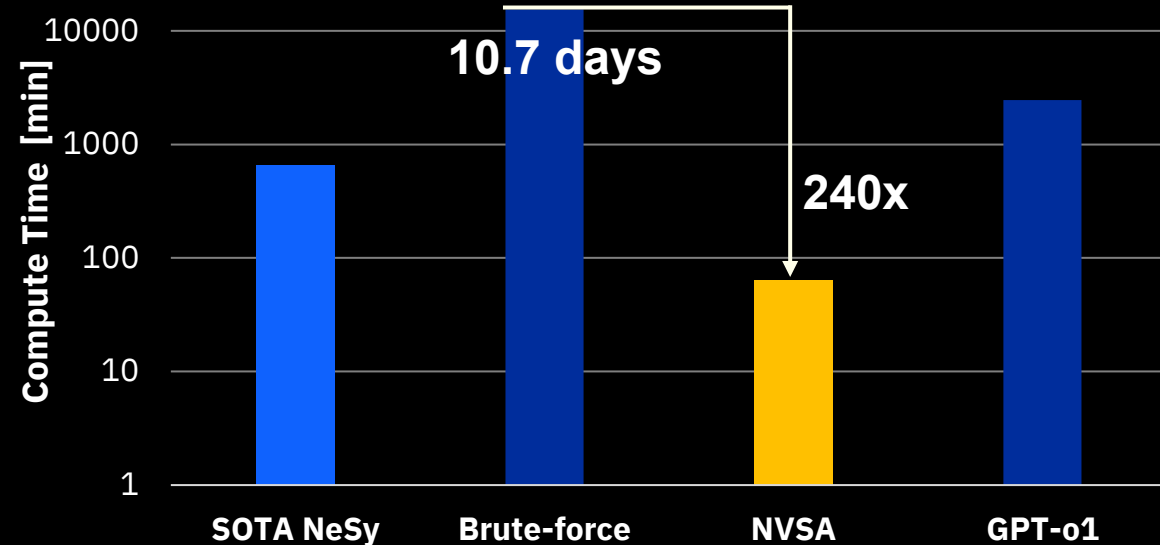
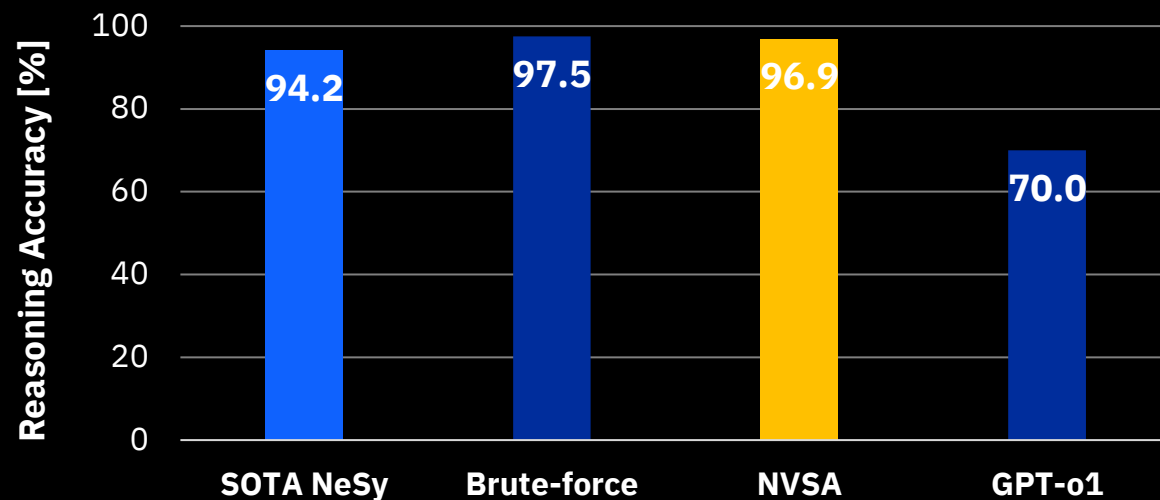
[Hersche et al., A Neuro-vector-symbolic architecture for solving Raven's progressive matrices, *Nature Machine Intelligence*, 2023]
Spotlighted in Quanta Magazine as one of the three biggest achievements of the 2023 in computer science

NVSA's reasoning is **240x** faster
enabling real-time inference on CPUs



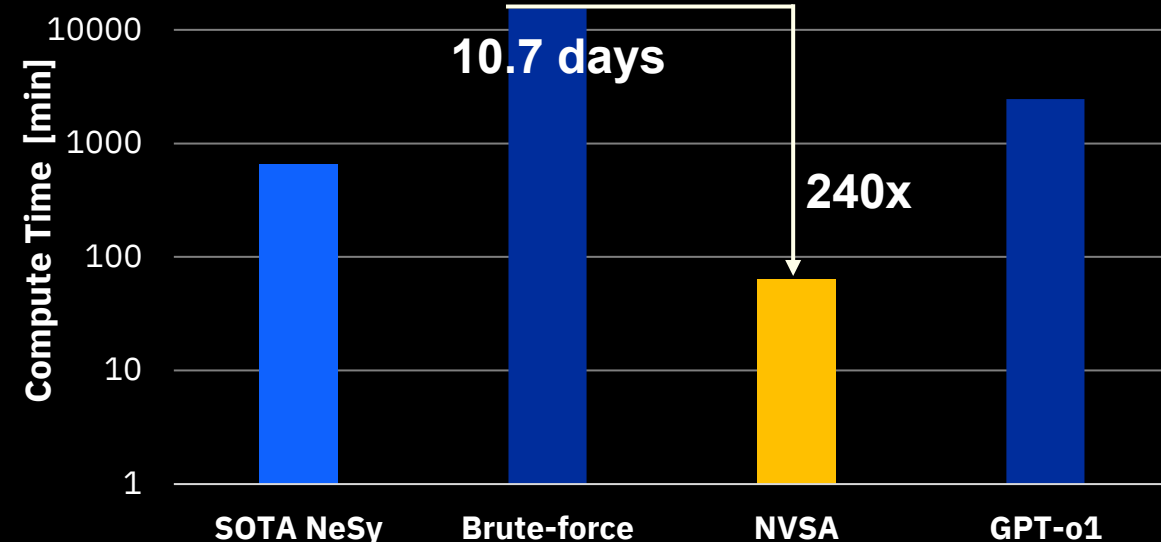
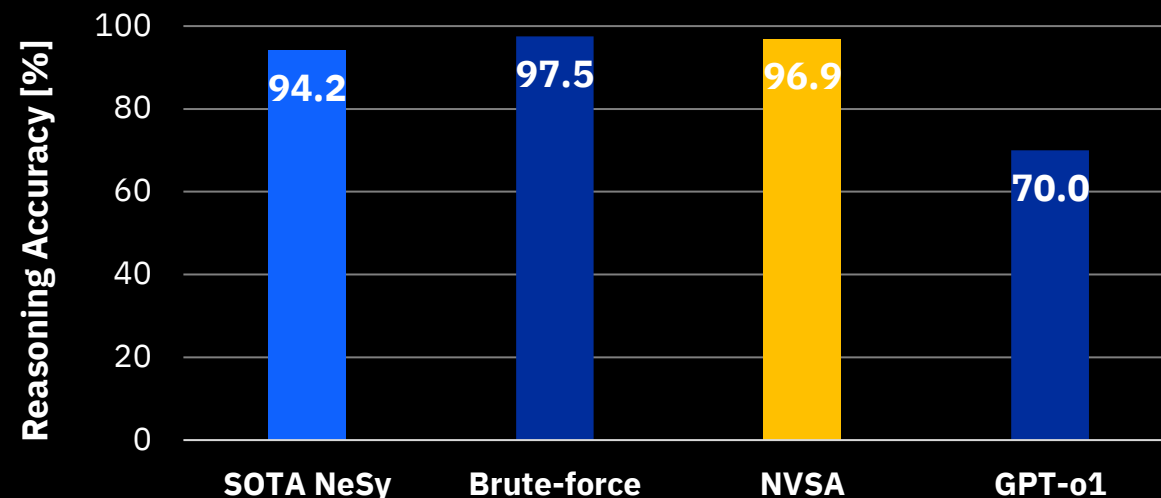
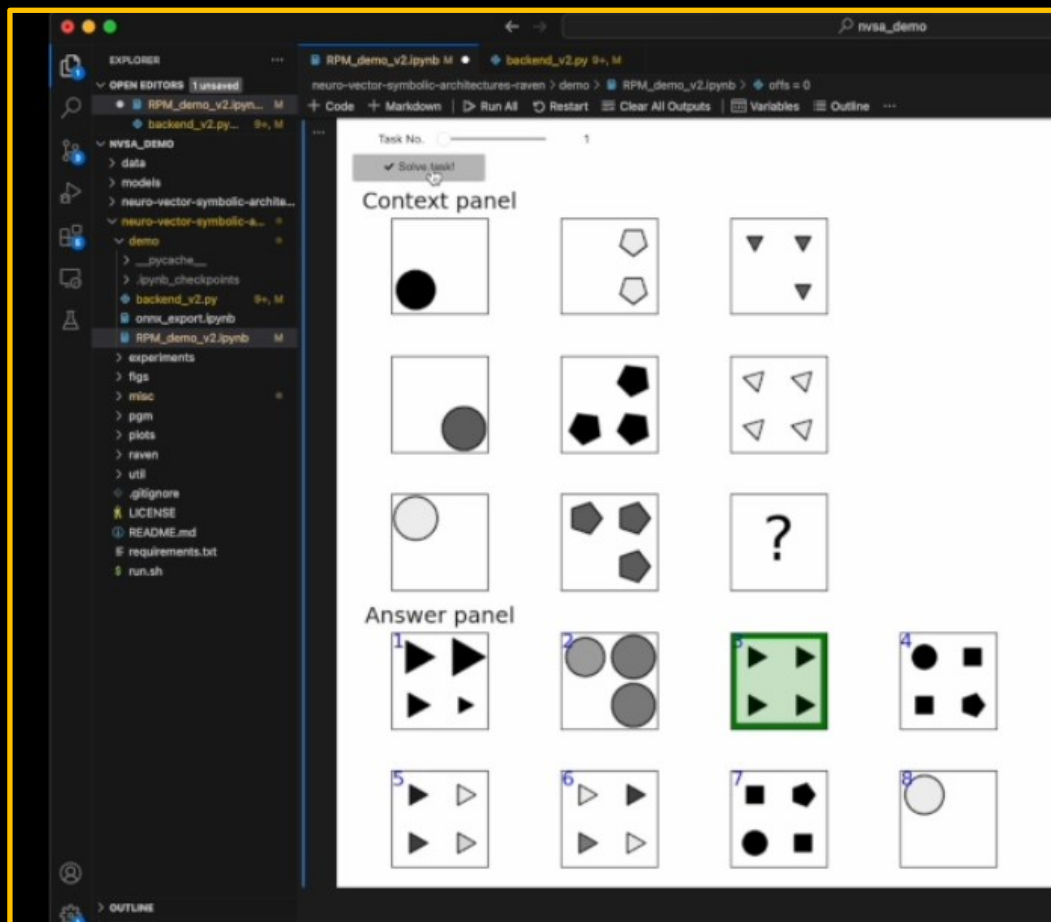
[Hersche et al., A Neuro-vector-symbolic architecture for solving Raven's progressive matrices, *Nature Machine Intelligence*, 2023]
Spotlighted in Quanta Magazine as one of the three biggest achievements of the 2023 in computer science

NVSA's reasoning is **240x** faster
enabling real-time inference on CPUs



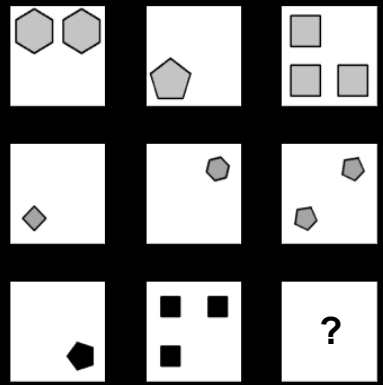
[Hersche et al., A Neuro-vector-symbolic architecture for solving Raven's progressive matrices, *Nature Machine Intelligence*, 2023]
Spotlighted in Quanta Magazine as one of the three biggest achievements of the 2023 in computer science

NVSA's reasoning is **240x** faster enabling real-time inference on CPUs



[Hersche et al., A Neuro-vector-symbolic architecture for solving Raven's progressive matrices, *Nature Machine Intelligence*, 2023]
Spotlighted in Quanta Magazine as one of the three biggest achievements of the 2023 in computer science

Toward learning to reason



Disentangling
neural
representations

NVSA (knowledge-based)

Mapping to
high-dimensional
space

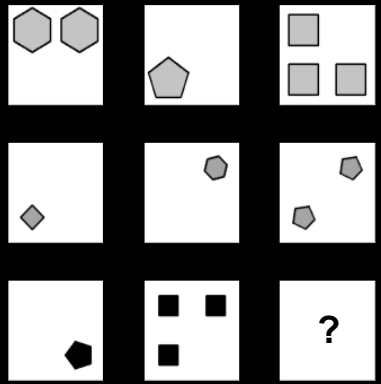


Rule
detection

Rule
execution

Select
answer
candidate

Toward learning to reason



Disentangling
neural
representations

NVSA (knowledge-based)

Mapping to
high-dimensional
space



Rule
detection

Rule
execution

Select
answer
candidate

$$\hat{a}_{3,3} = a_{3,1} = a_{3,2}$$

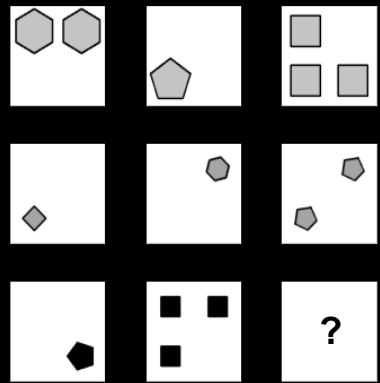
$$\hat{a}_{3,3} = a_{3,1} \otimes a_{3,2}$$

$$\hat{a}_{3,3} = a_{3,2} \otimes (a_{3,2} \oslash a_{3,1})$$

$$\hat{a}_{3,3} = (a_{1,2} \otimes a_{1,2} \otimes a_{1,3}) \oslash (a_{3,1} \otimes a_{3,2})$$

Hard-coded!

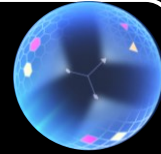
Toward learning to reason



Disentangling
neural
representations

NVSA (knowledge-based)

Mapping to
high-dimensional
space



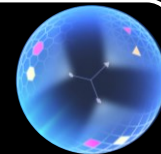
Rule
detection

Rule
execution

Select
answer
candidate

NVSA (L2R: learning-to-reason)

Mapping to
high-dimensional
space



Learning
to reason
 $\oplus \otimes \oslash$

Select
answer
candidate

$$\hat{a}_{3,3} = a_{3,1} = a_{3,2}$$

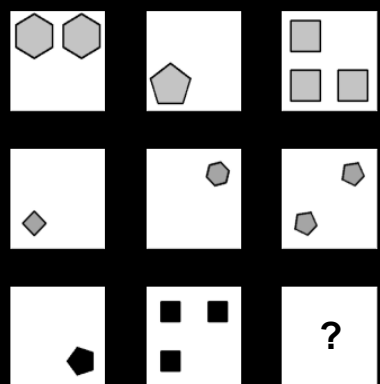
$$\hat{a}_{3,3} = a_{3,1} \otimes a_{3,2}$$

$$\hat{a}_{3,3} = a_{3,2} \otimes (a_{3,2} \oslash a_{3,1})$$

$$\hat{a}_{3,3} = (a_{1,2} \otimes a_{1,2} \otimes a_{1,3}) \oslash (a_{3,1} \otimes a_{3,2})$$

Hard-coded!

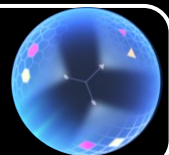
Toward learning to reason



Disentangling
neural
representations

NVSA (knowledge-based)

Mapping to
high-dimensional
space



Rule
detection

Rule
execution

Select
answer
candidate

$$\hat{a}_{3,3} = a_{3,1} = a_{3,2}$$

$$\hat{a}_{3,3} = a_{3,1} \otimes a_{3,2}$$

$$\hat{a}_{3,3} = a_{3,2} \otimes (a_{3,2} \oslash a_{3,1})$$

$$\hat{a}_{3,3} = (a_{1,2} \otimes a_{1,2} \otimes a_{1,3}) \oslash (a_{3,1} \otimes a_{3,2})$$

Hard-coded!

NVSA (L2R: learning-to-reason)

Mapping to
high-dimensional
space



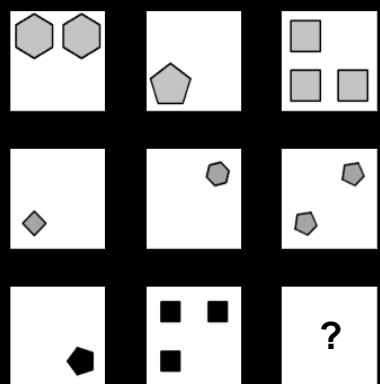
Learning
to reason
 $\oplus \otimes \oslash$

Select
answer
candidate

$$\prod_{k=1}^m \left(\sum_{i=1}^l w_k^i \cdot x_i + \sum_{j=1}^J u_k^j \cdot o_j + v_k e \right) \oslash \prod_{k=m+1}^n \left(\sum_{i=1}^l w_k^i \cdot x_i + \sum_{j=1}^J u_k^j \cdot o_j + v_k e \right)$$

Learning rules as an assignment problem

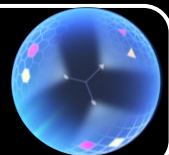
Toward learning to reason



Disentangling
neural
representations

NVSA (knowledge-based)

Mapping to
high-dimensional
space



Rule
detection

Rule
execution

Select
answer
candidate

$$\hat{a}_{3,3} = a_{3,1} = a_{3,2}$$

$$\hat{a}_{3,3} = a_{3,1} \otimes a_{3,2}$$

$$\hat{a}_{3,3} = a_{3,2} \otimes (a_{3,2} \oslash a_{3,1})$$

$$\hat{a}_{3,3} = (a_{1,2} \otimes a_{1,2} \otimes a_{1,3}) \oslash (a_{3,1} \otimes a_{3,2})$$

Hard-coded!

NVSA (L2R: learning-to-reason)

Mapping to
high-dimensional
space



Learning
to reason
 $\oplus \otimes \oslash$

Select
answer
candidate

$$\prod_{k=1}^m \left(\sum_{i=1}^l w_k^i \cdot x_i + \sum_{j=1}^J u_k^j \cdot o_j + v_k e \right) \oslash \prod_{k=m+1}^n \left(\sum_{i=1}^l w_k^i \cdot x_i + \sum_{j=1}^J u_k^j \cdot o_j + v_k e \right)$$

Learning rules as an assignment problem

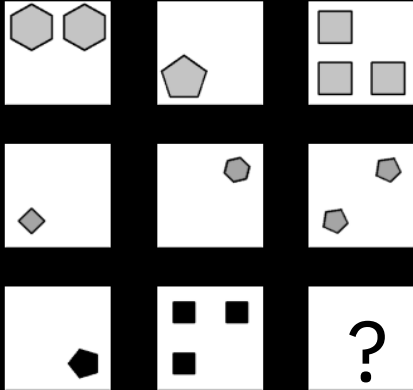
OOD tests:

GPT-3 87.6 ± 9.4

NVSA (L2R) 99.8 ± 0.4

Visual abstract reasoning:

3) OOD generalization



I-RAVEN (3x3)

System: Complete the Raven's progressive matrix:

User: Only return the missing numbers!

row 1: 2, 1, 3;

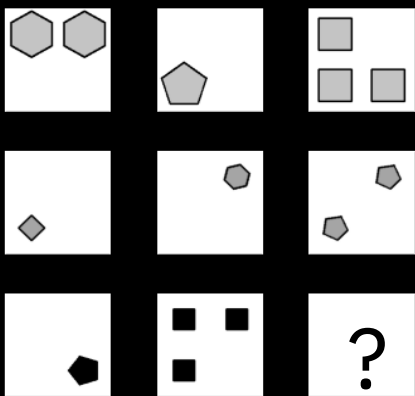
row 2: 1, 1, 2;

row 3: 1, 3,

Output: 4

Visual abstract reasoning:

3) OOD generalization



I-RAVEN (3x3)

System: Complete the Raven's progressive matrix:

User: Only return the missing numbers!

row 1: 2, 1, 3;

row 2: 1, 1, 2;

row 3: 1, 3,

Output: 4

Larger matrices
and inputs



I-RAVEN-X (3x10)

System: Complete the Raven's progressive matrix:

User: Only return the missing number!

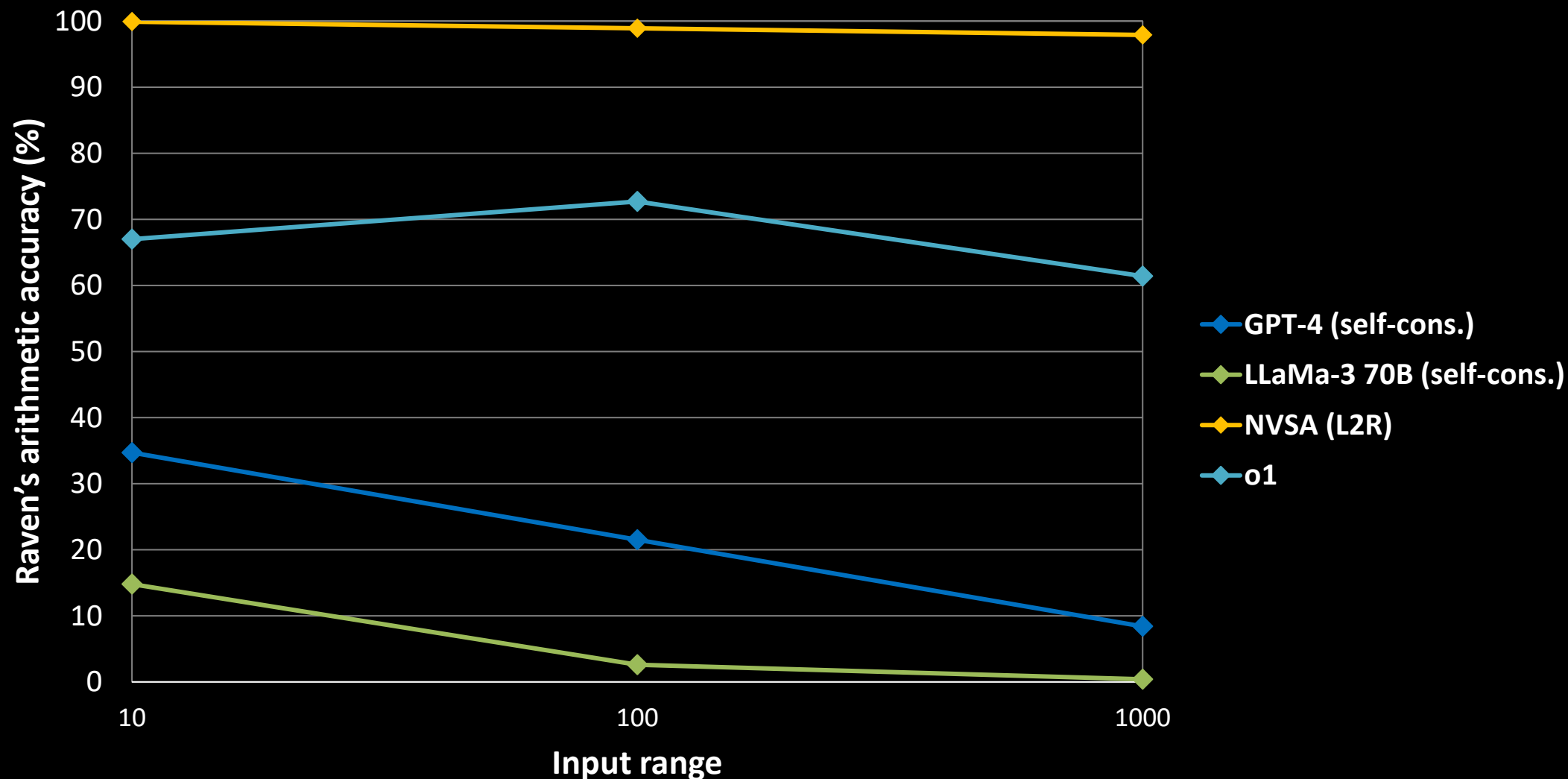
row 1: 769, 667, 0, 4, 2, 20, 63, 3, 5, 5;

row 2: 848, 0, 0, 0, 387, 2, 106, 7, 308, 38;

row 3: 611, 2, 0, 0, 0, 0, 0, 551, 0,

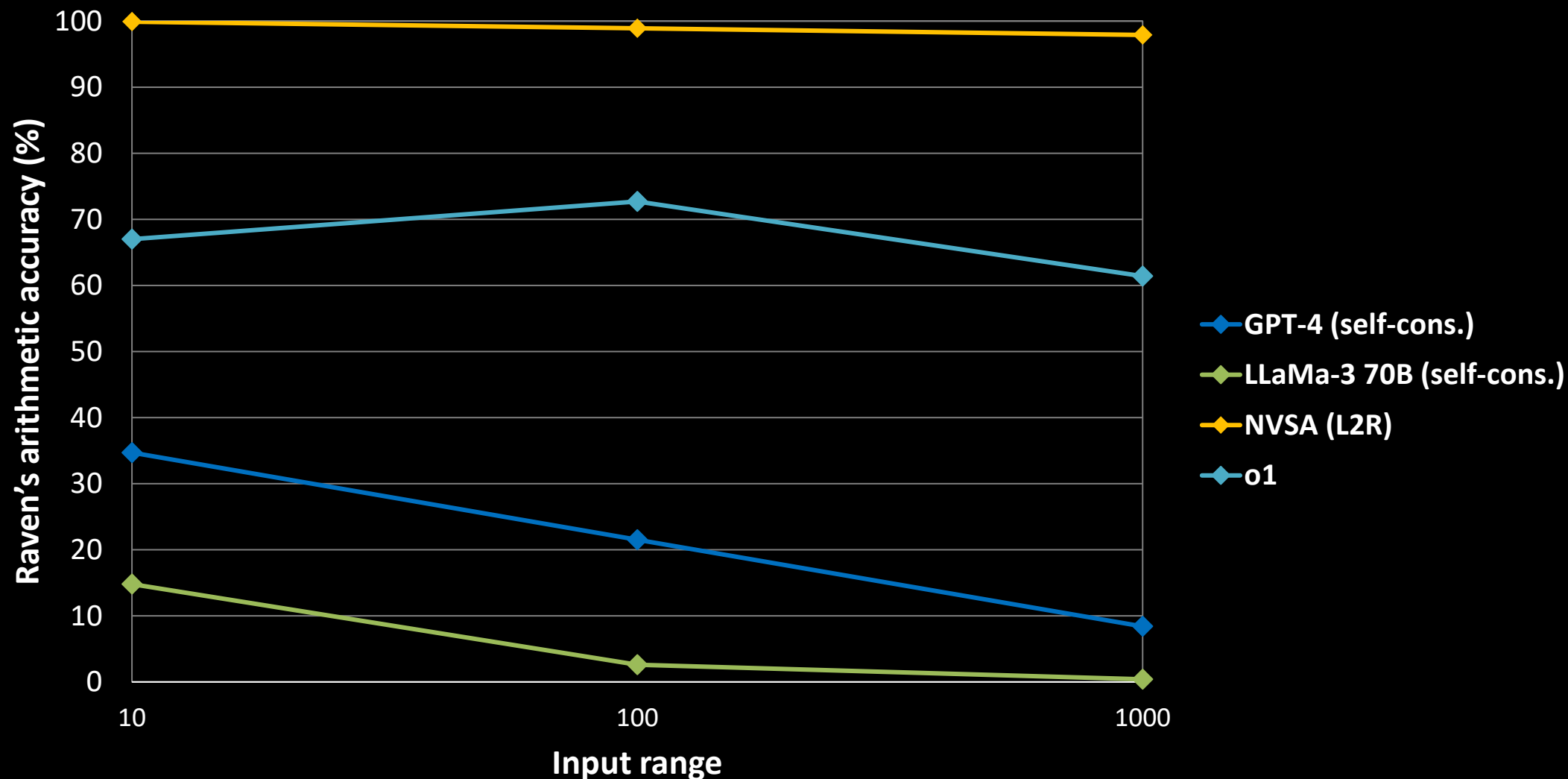
Output: 58

LLMs struggle with OOD

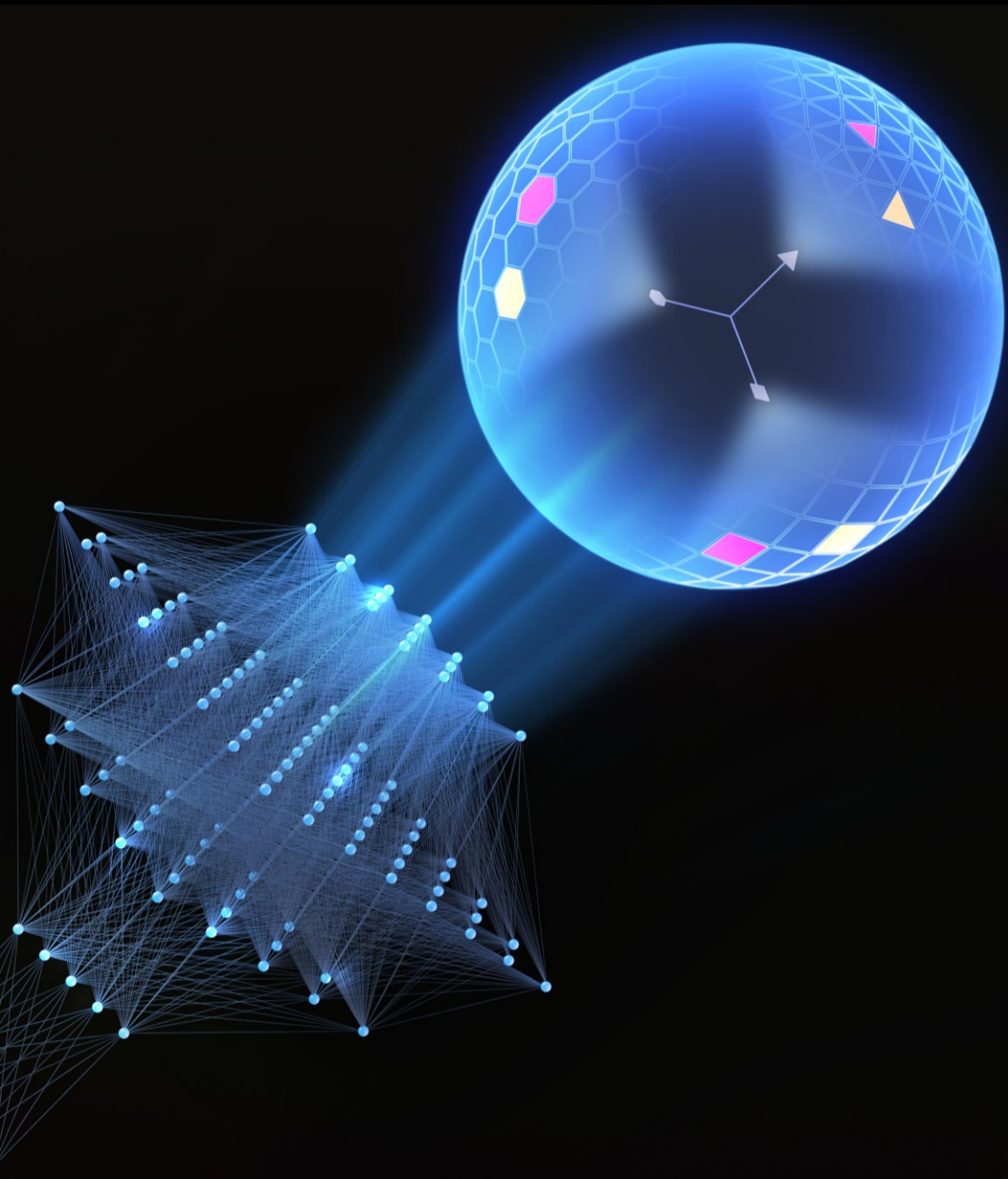


[Hersche et al., Towards learning to reason: comparing LLMs with neuro-symbolic on arithmetic relations in abstract reasoning, AAAI workshop Neural Reasoning and Mathematical Discovery, 2025]

LLMs struggle with OOD



[Hersche et al., Towards learning to reason: comparing LLMs with neuro-symbolic on arithmetic relations in abstract reasoning, AAAI workshop Neural Reasoning and Mathematical Discovery, 2025]



Summary

The roles of neuro-symbolic computing (albeit with distributed representations and operators) in machine intelligence:

- Enhancing **energy and compute efficiency in perception** by HW/SW co-design
- Fast and **scalable reasoning**
- Better **OOD generalization**