

Towards Cognitive AI Systems: Workload and Characterization of Neuro-Symbolic AI

Zishen Wan¹, Che-Kai Liu¹, Hanchen Yang¹, Ritik Raj¹, Chaojian Li¹, Haoran You¹, Yonggan Fu¹, Cheng Wan¹, Ananda Samajdar², Yingyan (Celine) Lin¹, Tushar Krishna¹, Arijit Raychowdhury¹

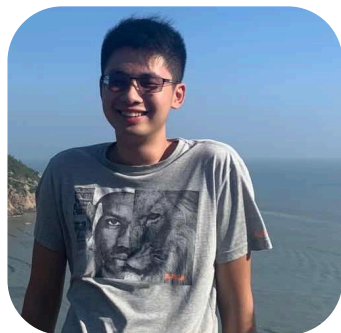
¹ Georgia Institute of Technology, GA ² IBM Research, NY

Email: zishenwan@gatech.edu

Our Team



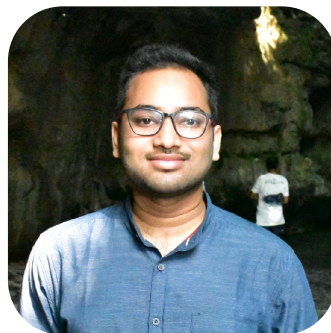
Zishen Wan



Che-Kai Liu



Hanchen Yang



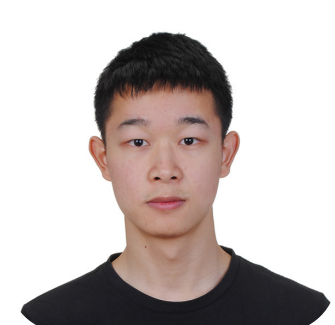
Ritik Raj



Chaojian Li



Haoran You



Yonggan Fu



Cheng Wan



Prof.
Celine Lin



Prof. Arijit
Raychowdhury



Prof. Tushar
Krishna



Ananda Samajdar



Neural Networks in Our Daily Life



Image Recognition



Speech Recognition



Language Translation



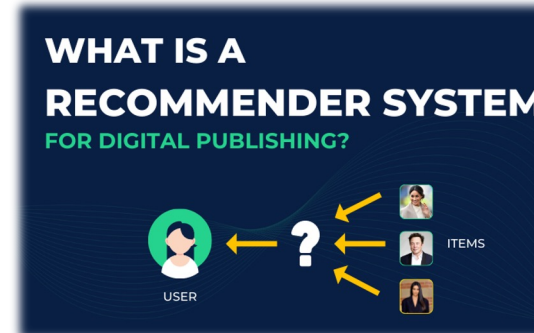
Autonomous Vehicle



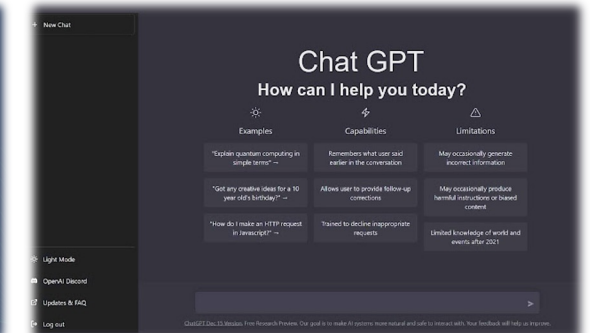
Medical Diagnosis



Financial Services

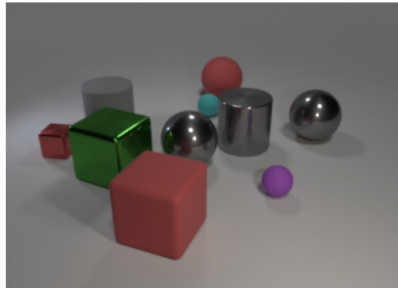


Recommendation Systems



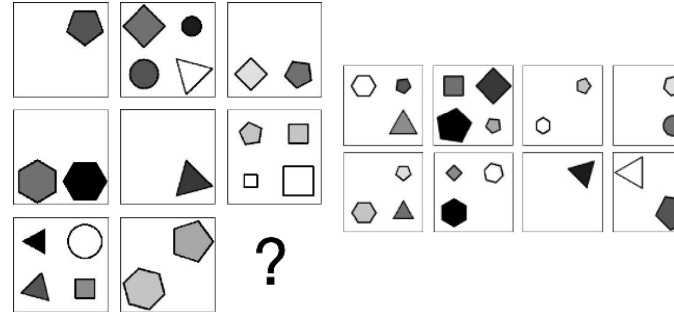
ChatGPT

But... Is That Enough?



(i) Remove all gray spheres. How many spheres are there? (3), (ii) Take away 3 cubes. How many objects are there? (7), (iii) How many blocks must be removed to get 1 block? (2)

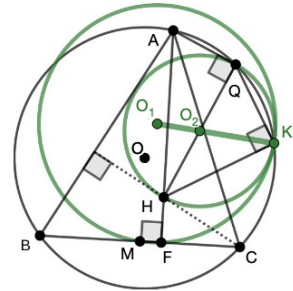
Complex Question Answering
NN accuracy: 50%



Abstract Reasoning
NN accuracy: 53%

IMO 2015 P3

“Let ABC be an acute triangle. Let (O) be its circumcircle, H its orthocenter, and F the foot of the altitude from A . Let M be the midpoint of BC . Let Q be the point on (O) such that $QH \perp QA$ and let K be the point on (O) such that $KH \perp KQ$. Prove that the circumcircles (O_1) and (O_2) of triangles FKM and KQH are tangent to each other.”



Automated Theorem Proving
NN accuracy: 0%



Interactive Learning
NN accuracy: 71%

Scenario
Imagine that a stranger will give Hank one thousand dollars to break all the windows in his neighbor's house without his neighbor's permission. Hank carries out the stranger's request.
Imagine that there are five people who are waiting in line to use a single-occupancy bathroom at a concert venue. Someone at the back of the line needs to throw up immediately. That person skips to the front of the line instead of waiting in the back.
At a summer camp, there is a pool. Right next to the pool is a tent where the kids at the camp have art class. The camp made a rule that there would be no cannonballing in the pool so that the art wouldn't get ruined by the splashing water. Today, there is a bee attacking this kid, and she needs to jump into the water quickly. This kid cannonballs into the pool.



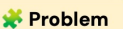
Ethical Decision Making
NN accuracy: 65%

Farmer John has N cows ($2 \leq N \leq 10^5$). Each cow has a breed that is either Guernsey or Holstein. As is often the case, the cows are standing in a line, numbered $1 \dots N$ in this order.

Over the course of the day, each cow writes down a list of cows. Specifically, cow i 's list contains the range of cows starting with herself (cow i) up to and including cow E_i ($i \leq E_i \leq N$).

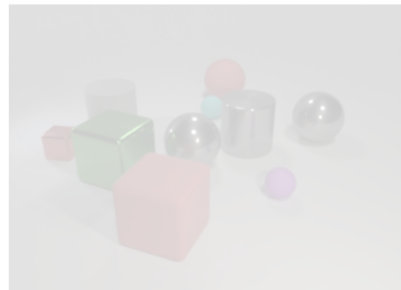
FJ has recently discovered that each breed of cow has exactly one distinct leader. FJ does not know who the leaders are, but he knows that each leader must have a list that includes all the cows of their breed, or the other breed's leader (or both).

Help FJ count the number of pairs of cows that could be leaders. It is guaranteed that there is at least one possible pair.

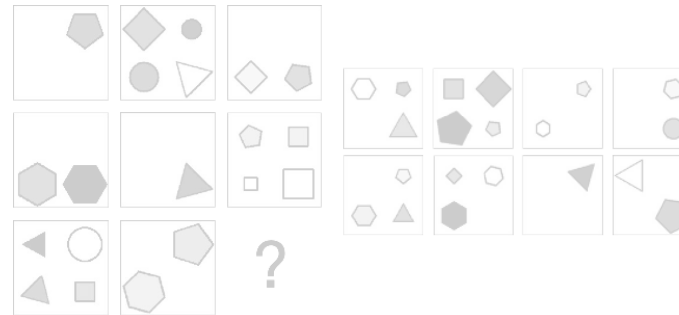


Competitive Programming
NN accuracy: 8.7%

But... Is That Enough?

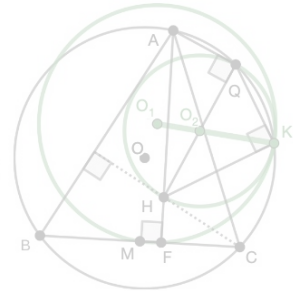


(i) Remove all gray spheres. How many spheres are there? (3), (ii) Take away 3 cubes. How many objects are there? (7), (iii) How many blocks must be removed to get 1 block? (2)



IMO 2015 P3

"Let ABC be an acute triangle. Let (O) be its circumcircle, H its orthocenter, and F the foot of the altitude from A . Let M be the midpoint of BC . Let Q be the point on (O) such that $QH \perp QA$ and let K be the point on (O) such that $KH \perp KQ$. Prove that the circumcircles (O_1) and (O_2) of triangles FKM and KQH are tangent to each other."



Complex Question Answering
NN accuracy: 50%

Abstract Reasoning
NN accuracy: 56%

Automated Theorem Proving
NN accuracy: 0%

Neuro-Symbolic AI



Interactive Learning
NN accuracy: 71%

Scenario
Imagine that a stranger will give Hank one thousand dollars to break all the windows in his neighbor's house without his neighbor's permission. Hank carries out the stranger's request.
Imagine that there are five people who are waiting in line to use a single-occupancy bathroom at a concert venue. Someone at the back of the line needs to throw up immediately. That person skips to the front of the line instead of waiting in the back.
At a summer camp, there is a pool. Right next to the pool is a tent where the kids at the camp have art class. The camp made a rule that there would be no cannonballing in the pool so that the art wouldn't get ruined by the splashing water. Today, there is a bee attacking this kid, and she needs to jump into the water quickly. This kid cannonballs into the pool.



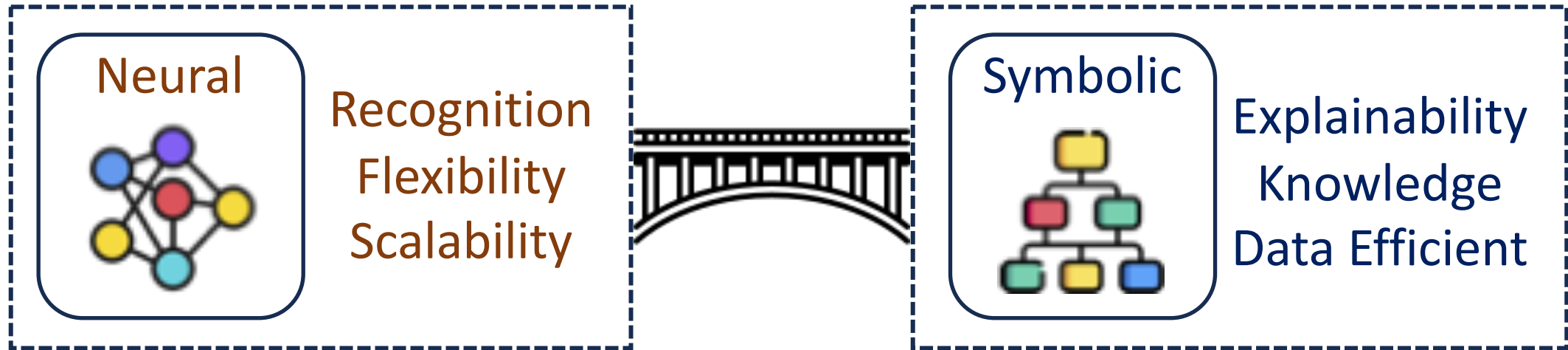
Ethical Decision Making
NN accuracy: 65%

Farmer John has N cows ($2 \leq N \leq 10^5$). Each cow has a breed that is either Guernsey or Holstein. As is often the case, the cows are standing in a line, numbered $1 \dots N$ in this order.
Over the course of the day, each cow writes down a list of cows. Specifically, cow i 's list contains the range of cows starting with herself (cow i) up to and including cow E_i ($i \leq E_i \leq N$).
FJ has recently discovered that each breed of cow has exactly one distinct leader. FJ does not know who the leaders are, but he knows that each leader must have a list that includes all the cows of their breed, or the other breed's leader (or both).
Help FJ count the number of pairs of cows that could be leaders. It is guaranteed that there is at least one possible pair.

Problem

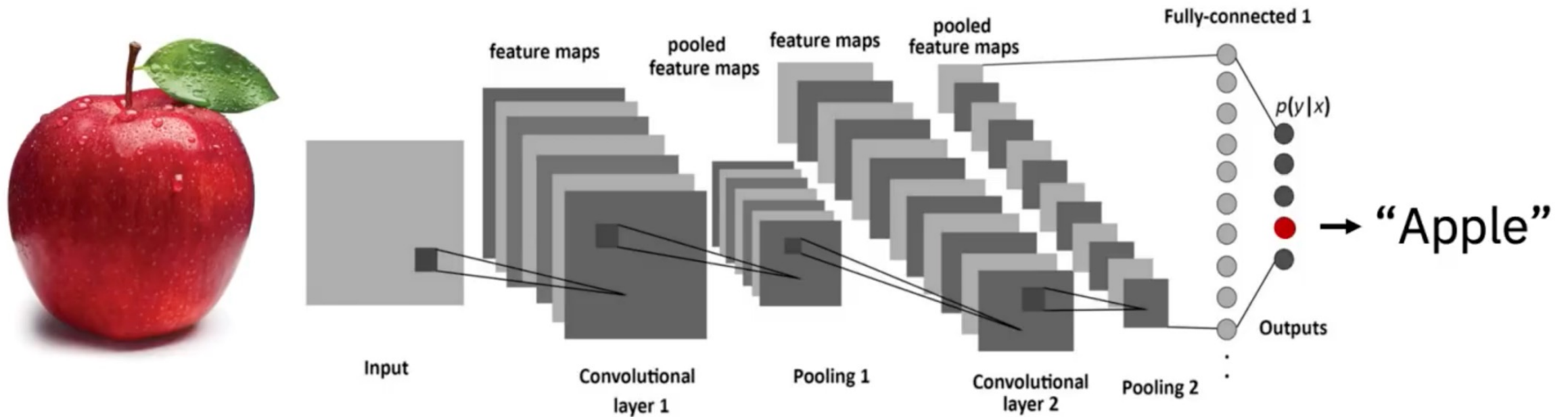
Competitive Programming
NN accuracy: 8.7%

What is Neuro-Symbolic AI?



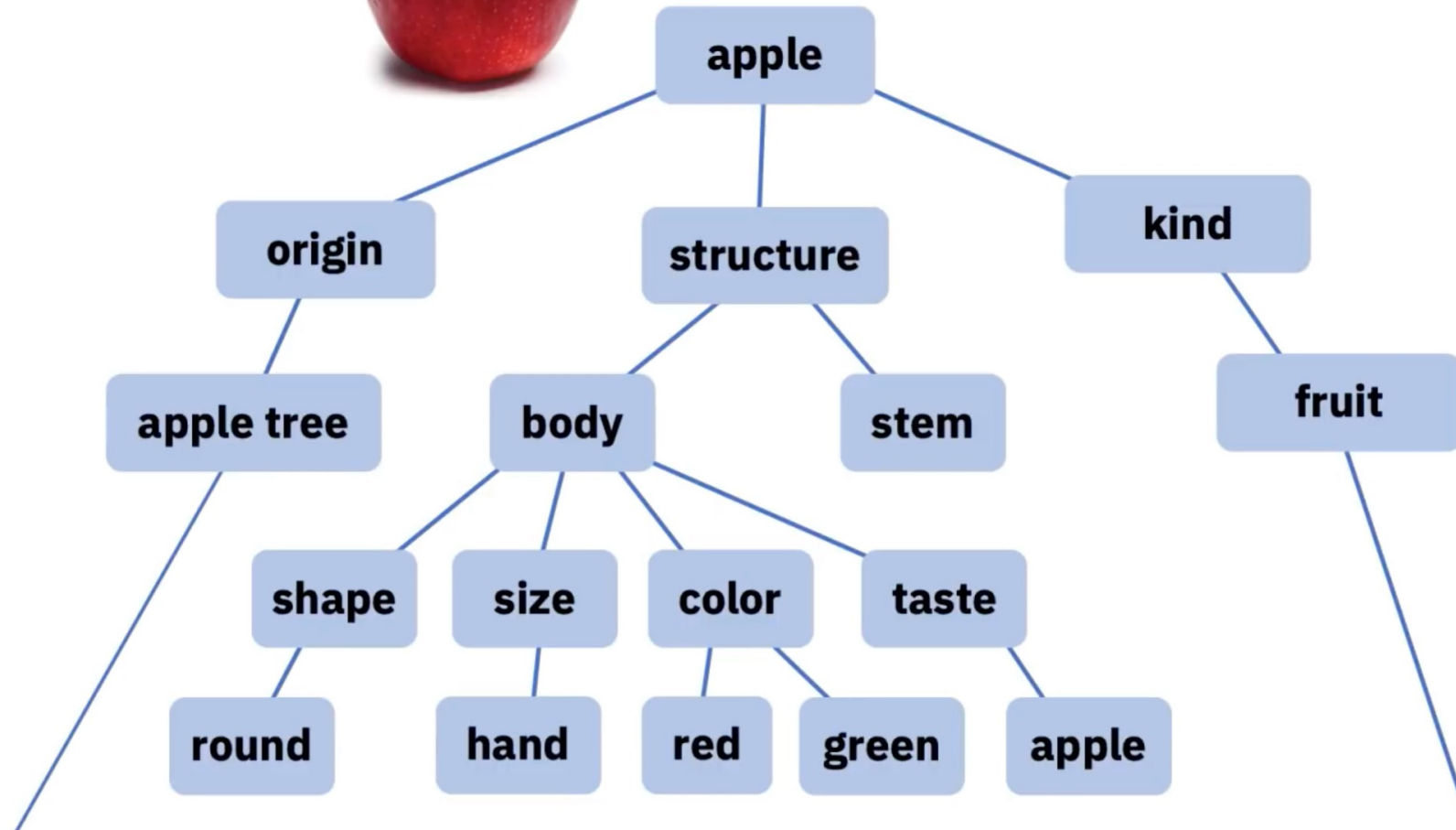
Towards Cognitive and Trustworthy AI Systems

Neural Network



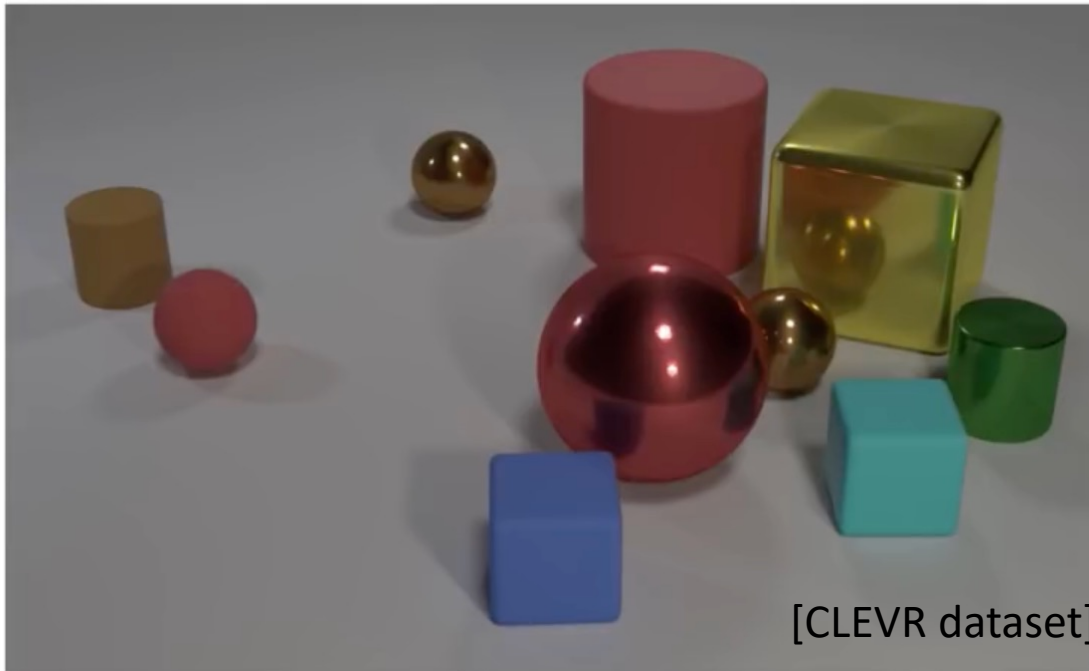
Slide Adapted from MIT 6.S191: Neurosymbolic AI

Symbolic AI



Slide Adapted from MIT 6.S191: Neurosymbolic AI

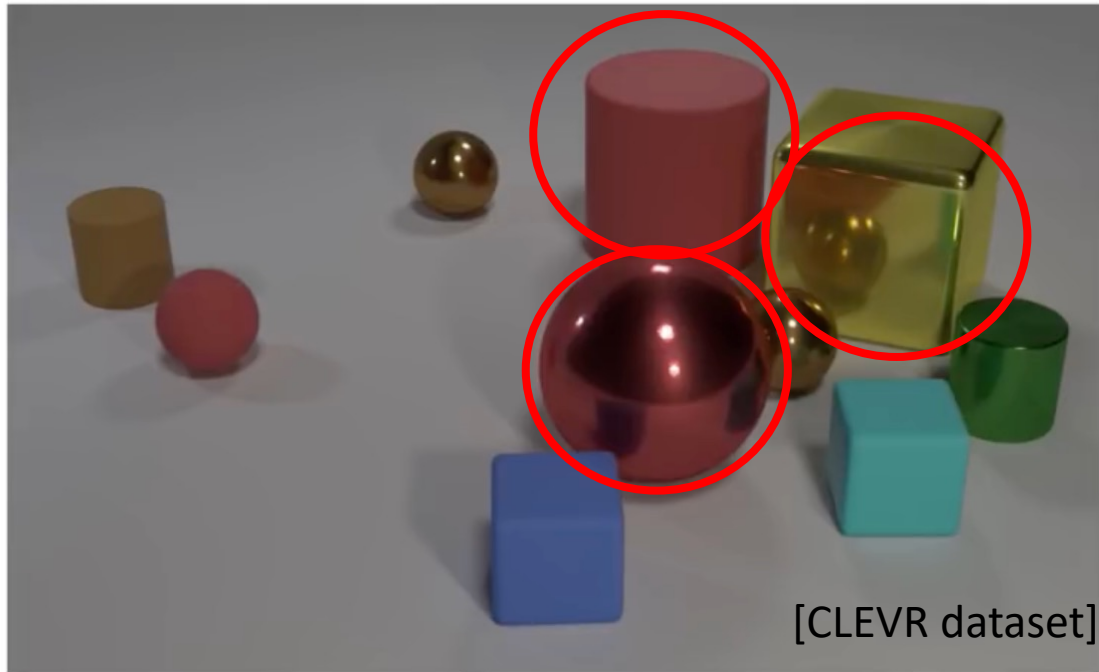
Neuro-Symbolic AI Example: Visual Reasoning



Question: *Are there an equal number of large things and metal spheres?*

Slide Adapted from MIT 6.S191: Neurosymbolic AI

Neuro-Symbolic AI Example: Visual Reasoning



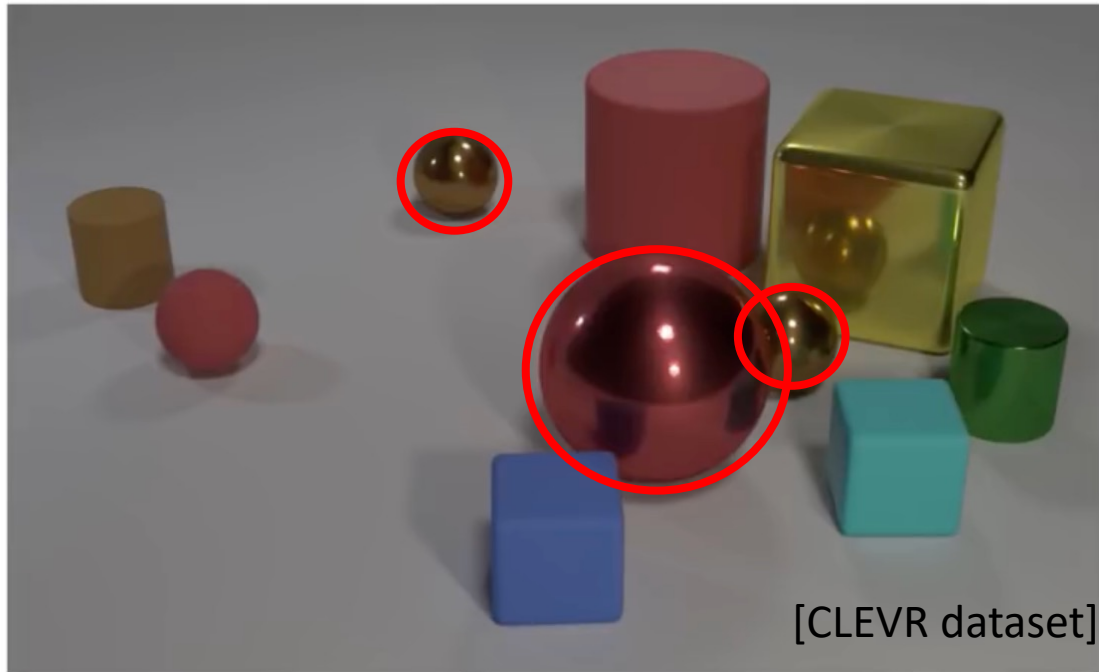
Question: *Are there an equal number of large things and metal spheres?*

3 large things!



Slide Adapted from MIT 6.S191: Neurosymbolic AI

Neuro-Symbolic AI Example: Visual Reasoning



Question: *Are there an equal number of large things and metal spheres?*

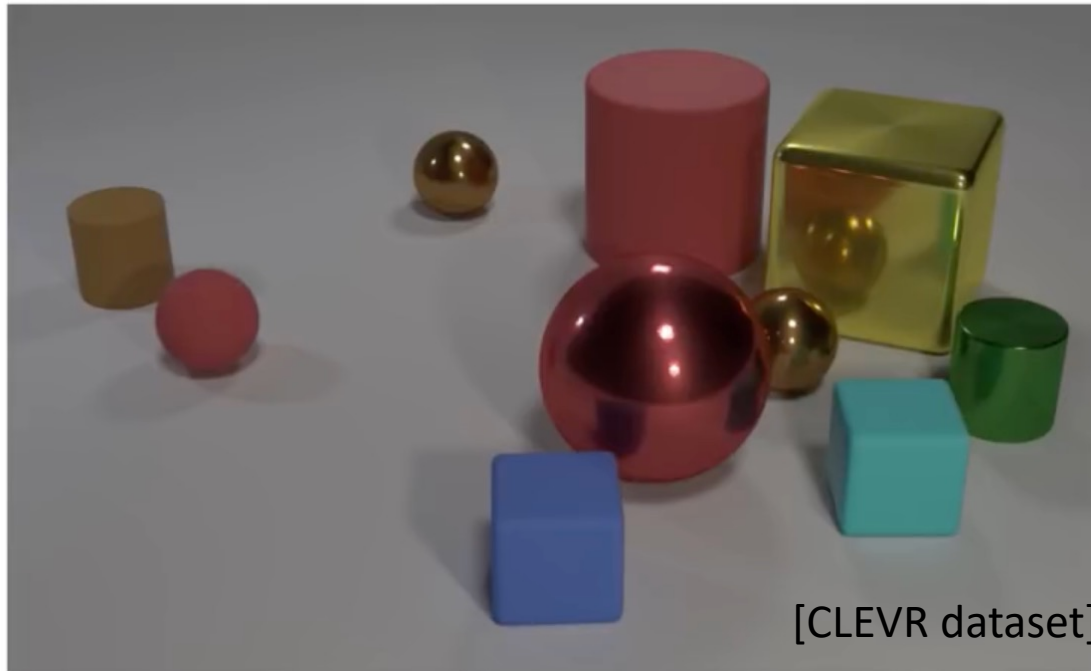
3 large things!

3 metal spheres!

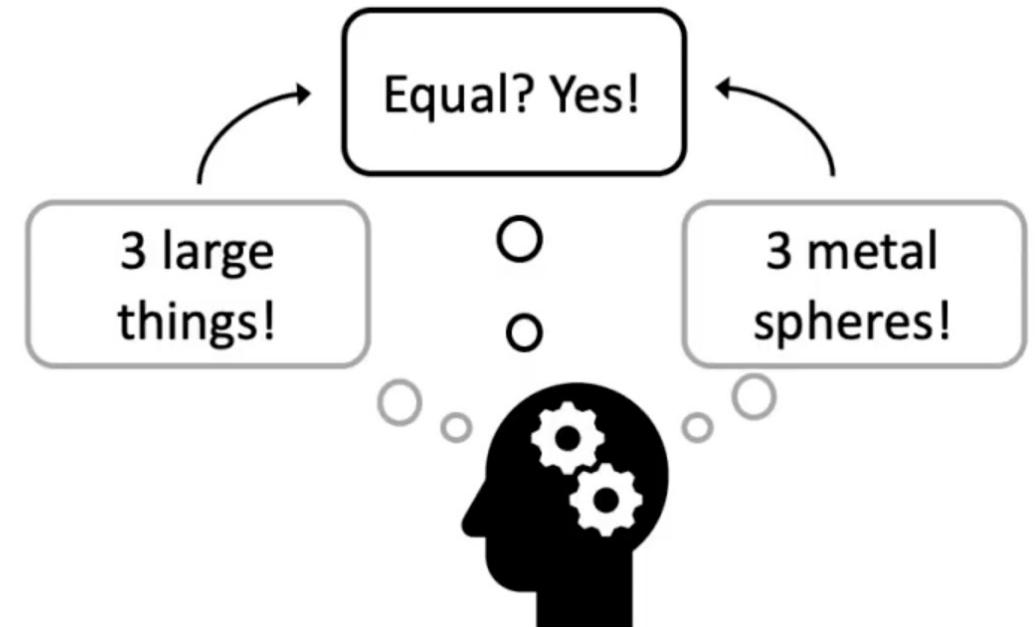


Slide Adapted from MIT 6.S191: Neurosymbolic AI

Neuro-Symbolic AI Example: Visual Reasoning

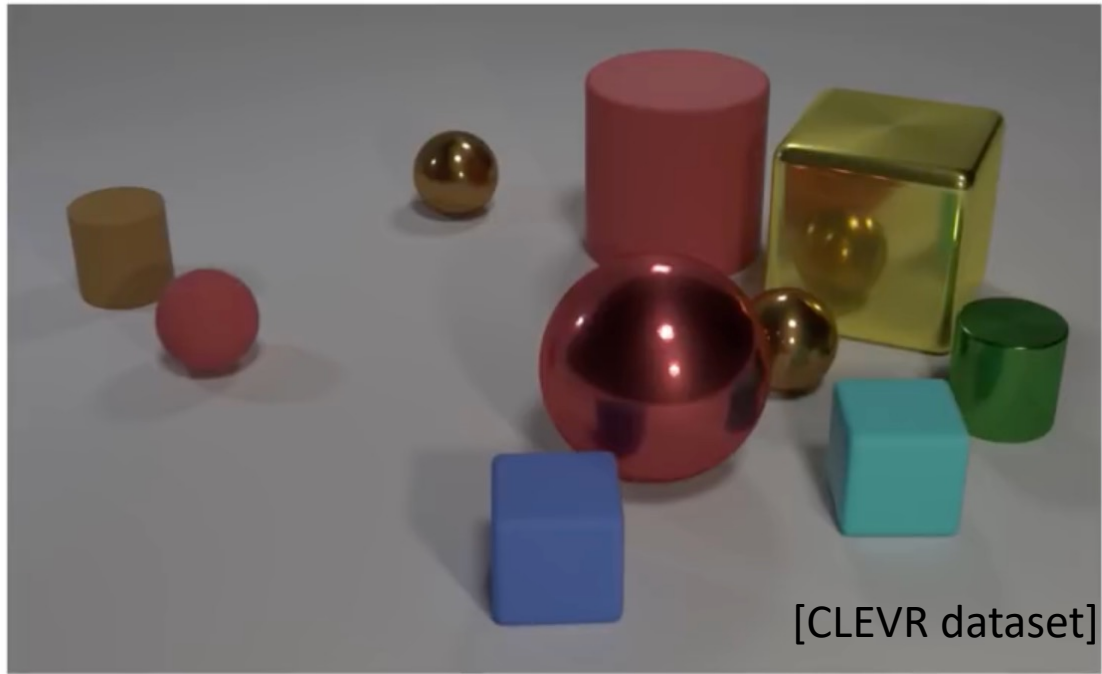


Question: *Are there an equal number of large things and metal spheres?*



Slide Adapted from MIT 6.S191: Neurosymbolic AI

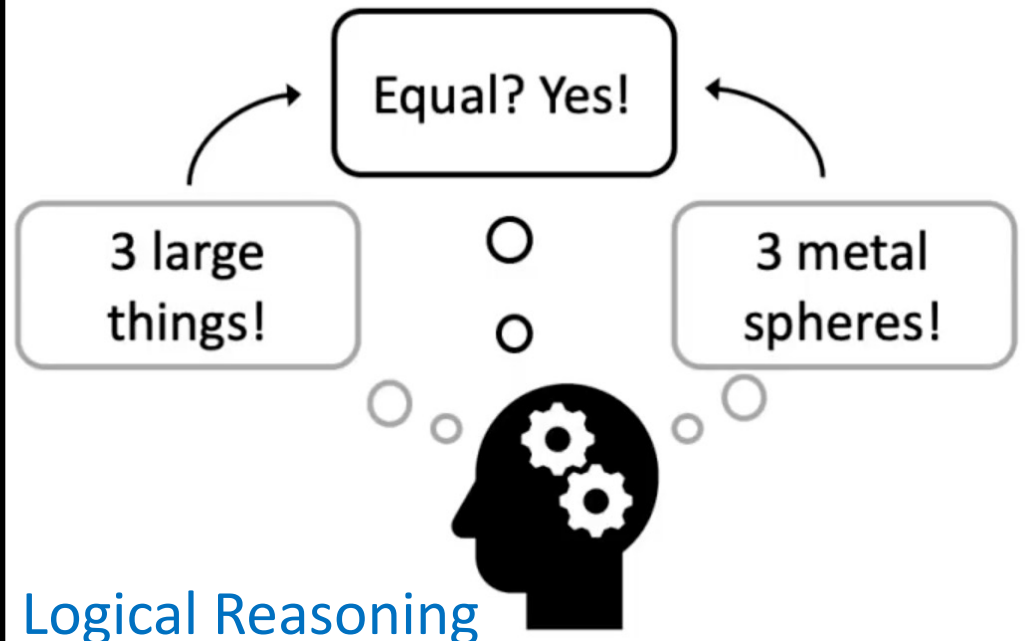
Neuro-Symbolic AI Example: Visual Reasoning



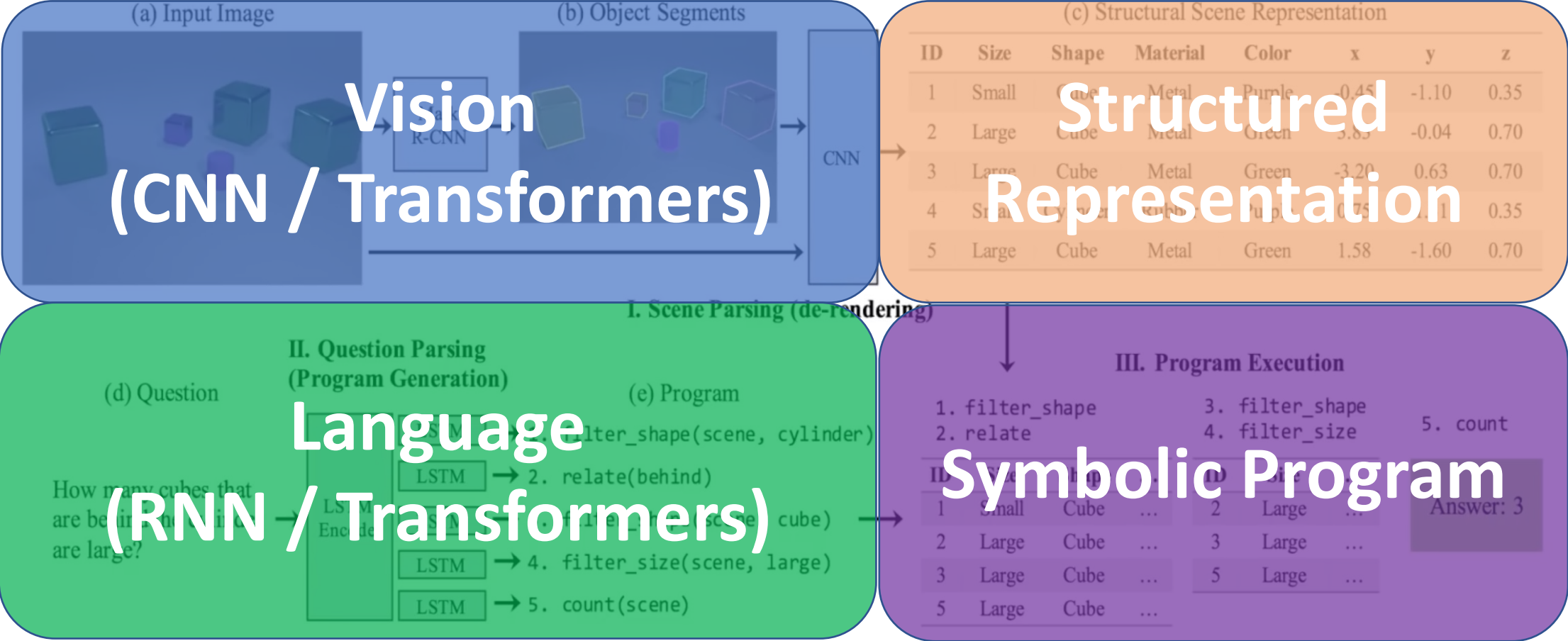
Visual Perception

Question Understanding

Question: *Are there an equal number of large things and metal spheres?*



Neuro-Symbolic AI Example: Visual Reasoning

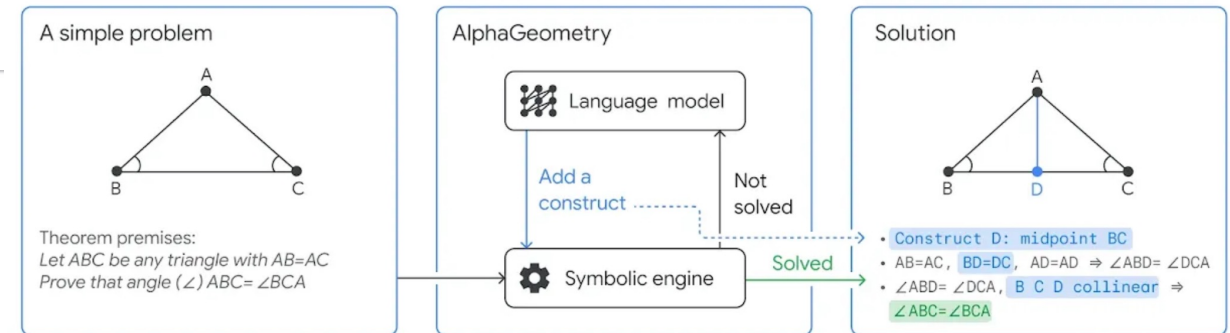


Other Examples

AlphaGeometry: An Olympiad-level AI system for geometry

17 JANUARY 2024
Trieu Trinh and Thang Luong

Share



LLM: construct generation
Symbolic: deductive reasoning

Eval on 30 Int. Math Olympics (IMO) problems:

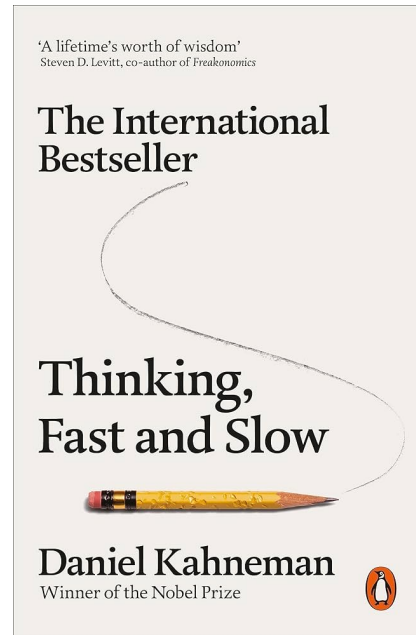
- **GPT-4:** 0/30
- **AlphaGeometry (Neuro-Symbolic):** 25/30
- **Human Gold Medalist:** 26/30

Trinh et al, "Solving Olympiad Geometry without Human Demonstrations", Nature 2024

Relationship to Human Minds



**Daniel Kahneman
(1934-2024)**



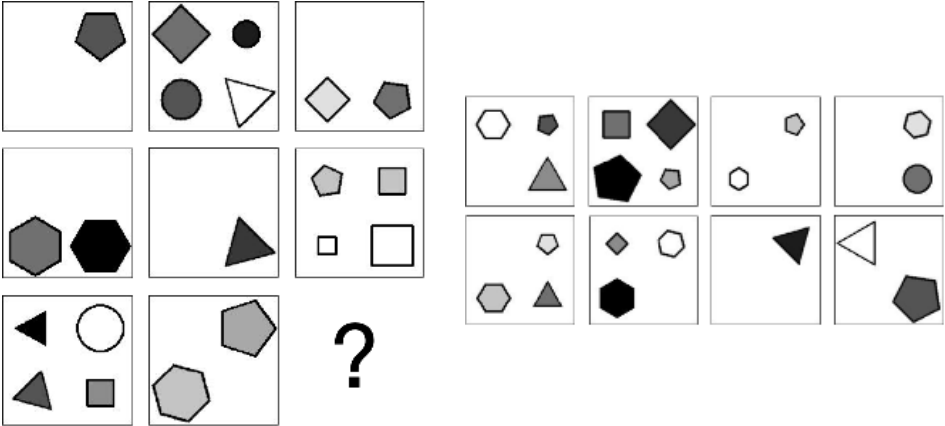
AlphaGeometry adopts a neuro-symbolic approach

AlphaGeometry is a neuro-symbolic system made up of a neural language model and a symbolic deduction engine, which work together to find proofs for complex geometry theorems. Akin to the idea of "[thinking, fast and slow](#)", one system provides fast, "intuitive" ideas, and the other, more deliberate, rational decision-making.

However...

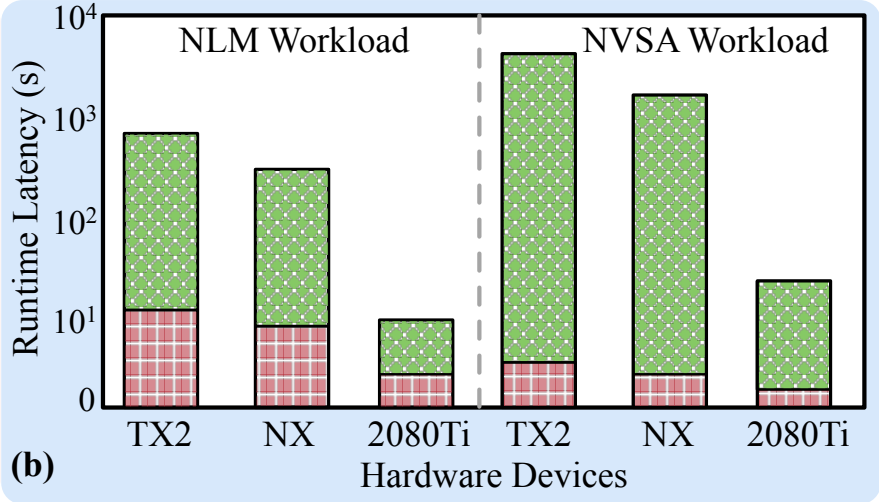


These neuro-symbolic approaches are typically very slow



Spatial-Temporal Abstract Reasoning

- ResNet accuracy: 53%
- GPT-4 accuracy: 84%
- Neuro-Symbolic accuracy: 98%

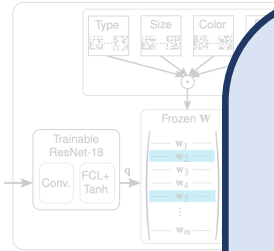


The neuro-symbolic approach takes ~100s even on desktop GPU, ~700s on Jetson TX2

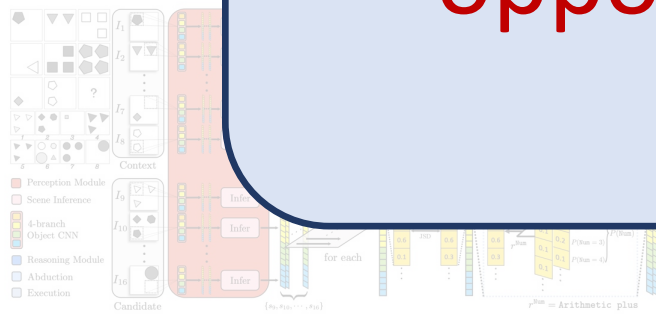
Lots of Neuro-Symbolic Algorithms

What's the **system behavior** and **co-design opportunities** of **Neuro-Symbolic AI**?

(b) NVSA frontend: perception



Neuro-Symbolic

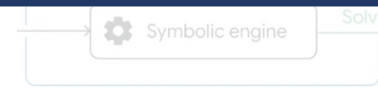


Probabilistic Abduction^[5]

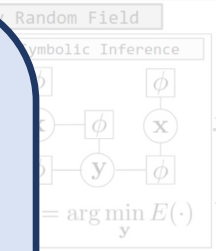
(A) Which ϕ is $\phi(\text{computer} \rightarrow \text{closed}) \rightarrow \text{closed}$

Image Translation via VSA^[6]

$\mathcal{G}(x)$



AlphaGeometry^[7]



Soft Logic^[4]



Logic-LM^[8]

[1] Hersche et al, Nature MI 2023; [2] Hoang et al, AAAI 2022; [3] Badreddine et al, AI 2022; [4] Pryor et al, IJCAI 2023

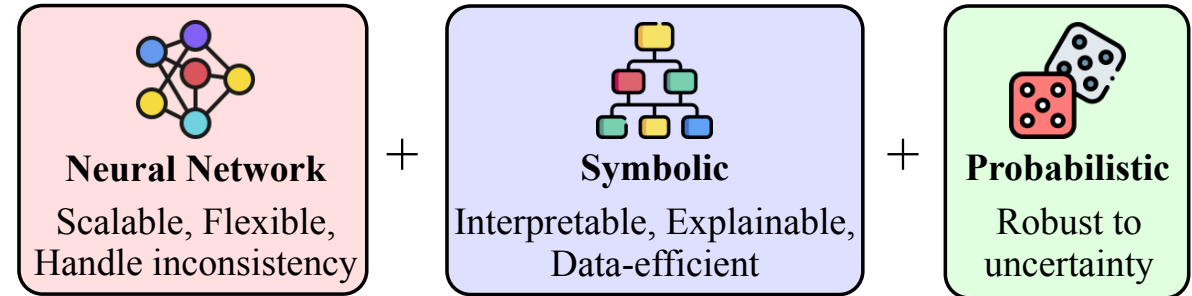
[5] Zhang et al, CVPR 2021; [6] Theiss et al, ECCV 2023; [7] Trinh et al, Nature 2024; [8] Pan et al, EMNLP 2023

Objective of this Work

Workload and Characterization of Neuro-Symbolic AI

Categorize Neuro-Symbolic Algorithms

Neuro-Symbolic AI Algorithms

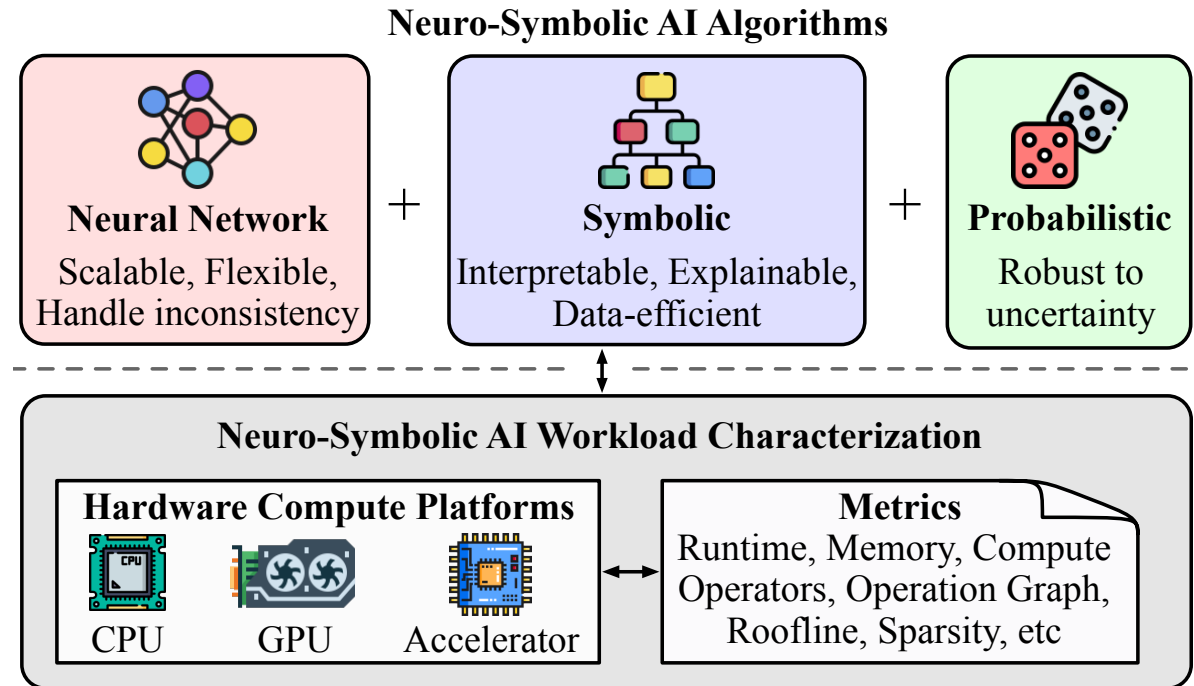


Objective of this Work

Workload and Characterization of Neuro-Symbolic AI

Categorize Neuro-Symbolic Algorithms

Understand Computational Behavior of Neuro-Symbolic Workloads



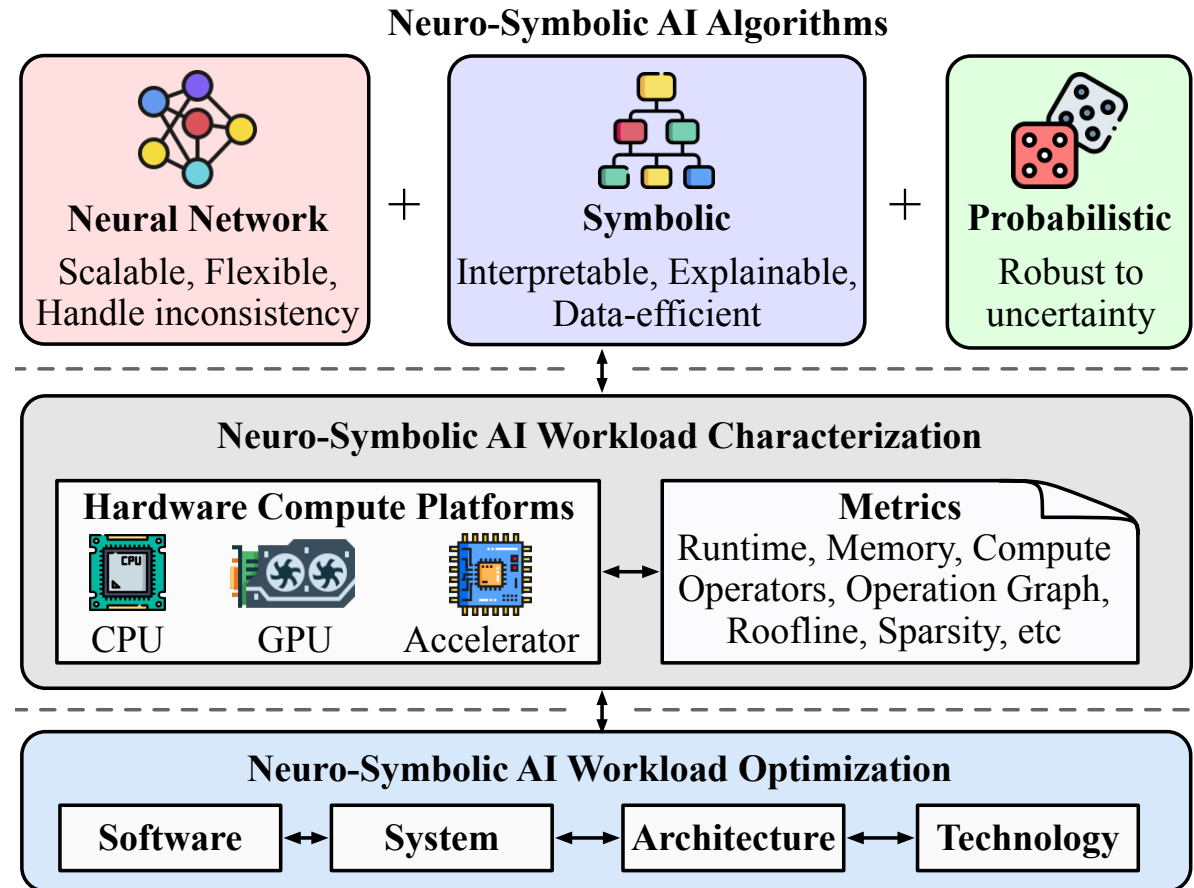
Objective of this Work

Workload and Characterization of Neuro-Symbolic AI

Categorize Neuro-Symbolic Algorithms

Understand Computational Behavior of Neuro-Symbolic Workloads

Identify Co-Design Opportunities



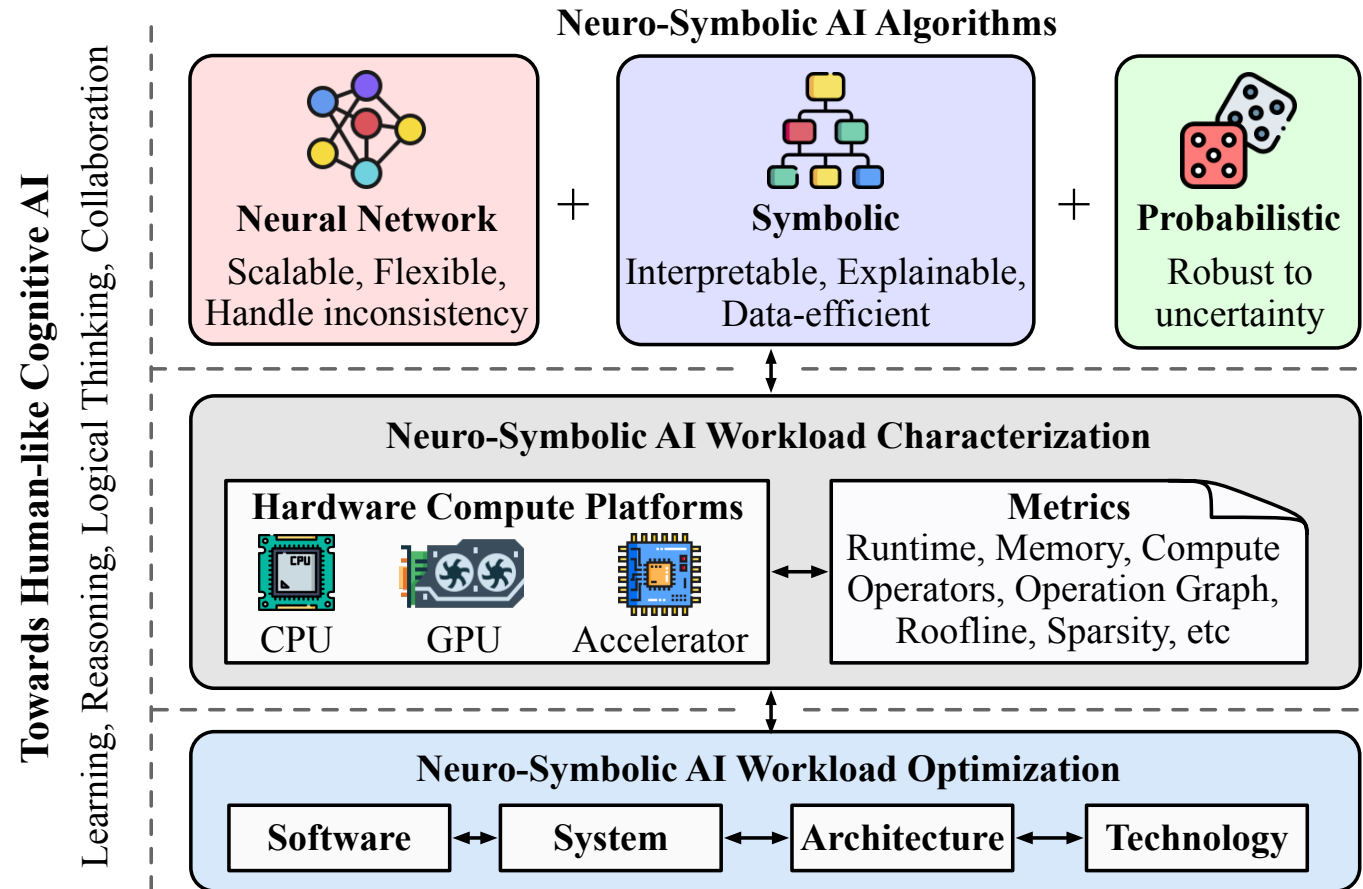
Objective of this Work

Workload and Characterization of Neuro-Symbolic AI

Categorize Neuro-Symbolic Algorithms

Understand Computational Behavior of Neuro-Symbolic Workloads

Identify Co-Design Opportunities



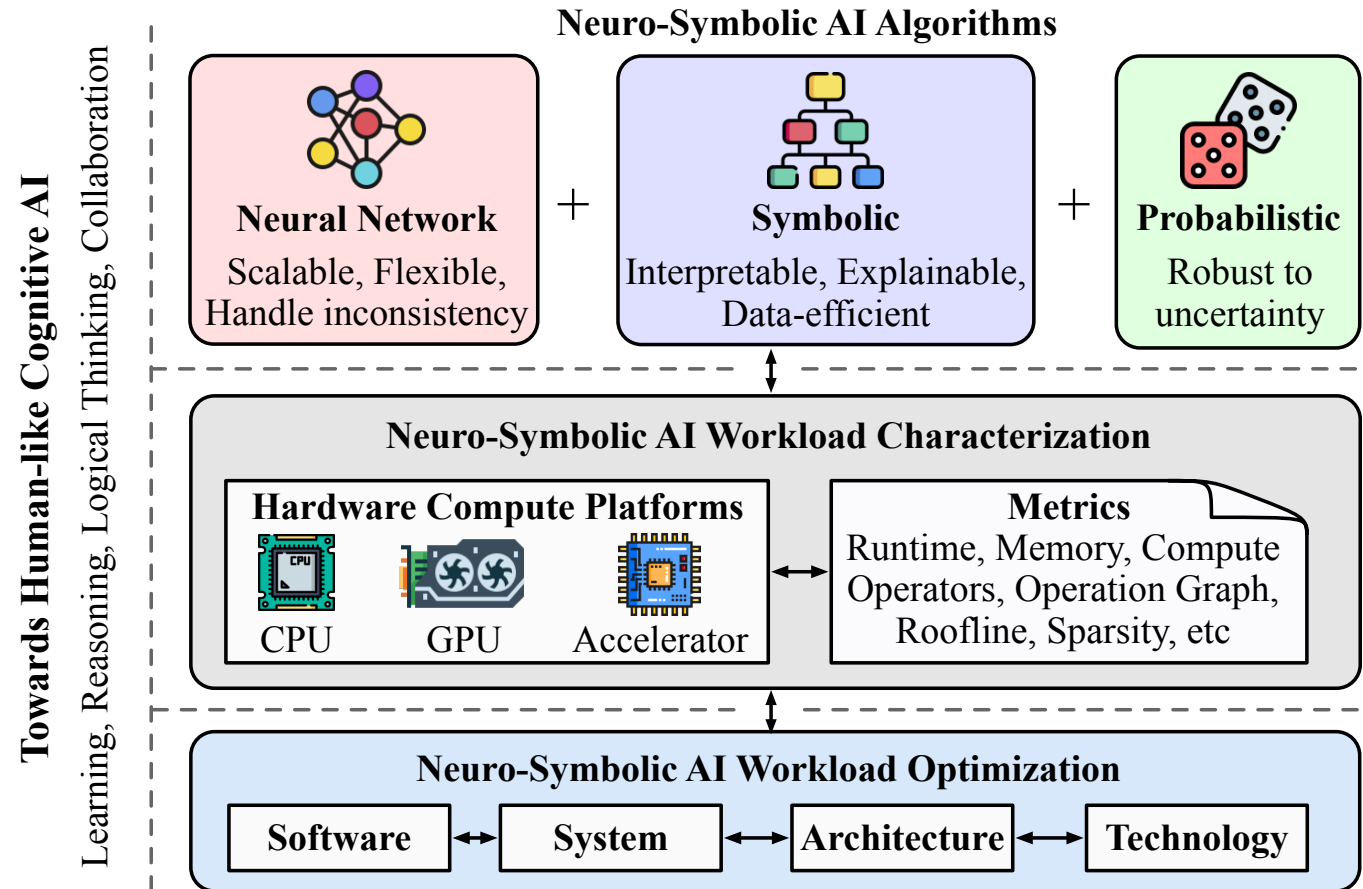
Objective of this Work

Workload and Characterization of Neuro-Symbolic AI

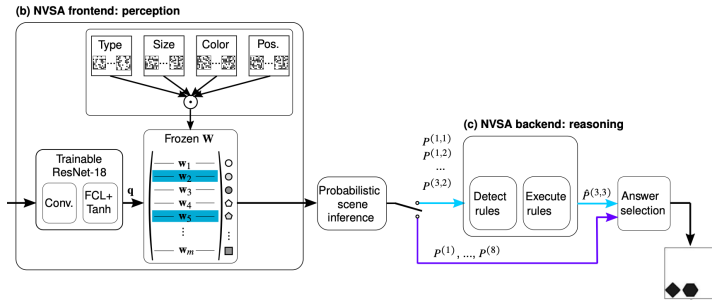
Categorize Neuro-Symbolic Algorithms

Understand Computational Behavior of Neuro-Symbolic Workloads

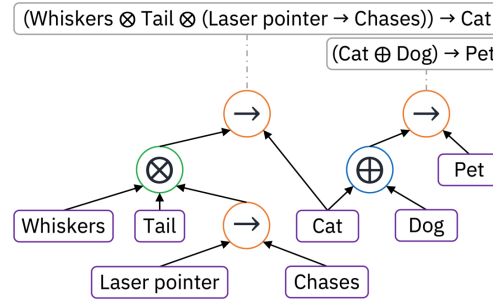
Identify Co-Design Opportunities



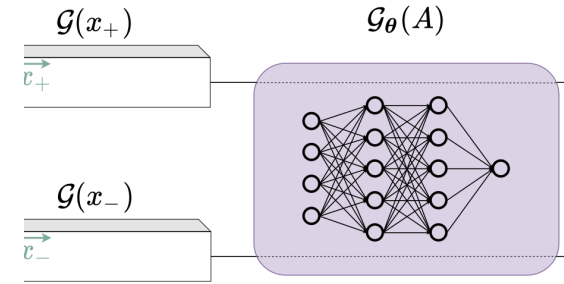
Lots of Neuro-Symbolic Algorithms



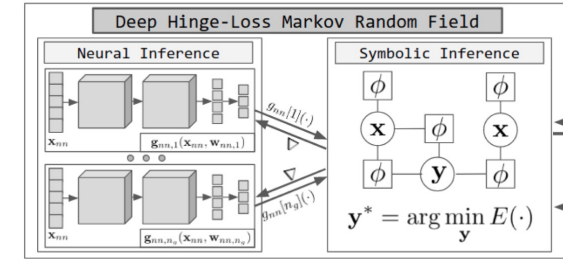
Neuro-Vector-Symbolic Arch



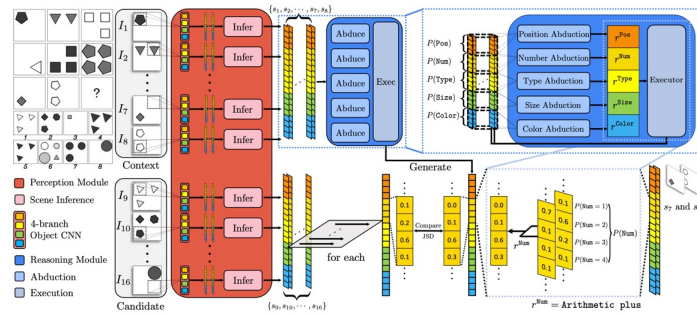
Logical Neural Network



Logical Tensor Network



Neural Probabilistic Soft Logic



Probabilistic Abduction

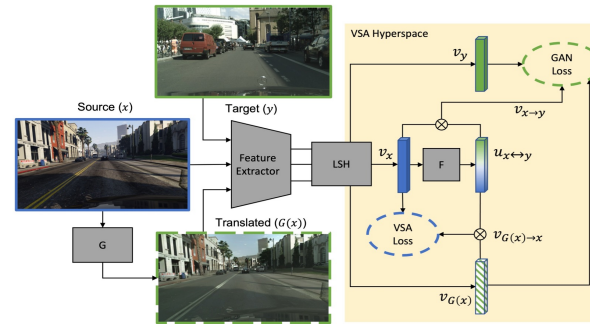
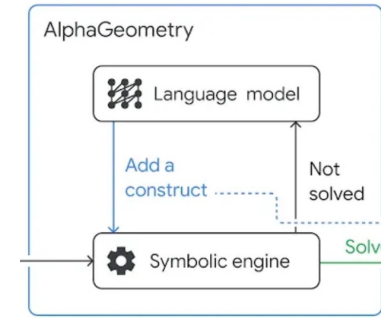
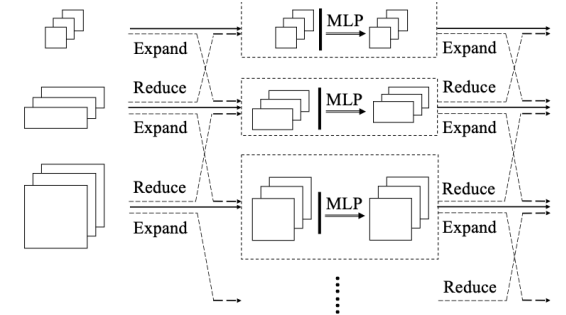


Image Translation via VSA



AlphaGeometry

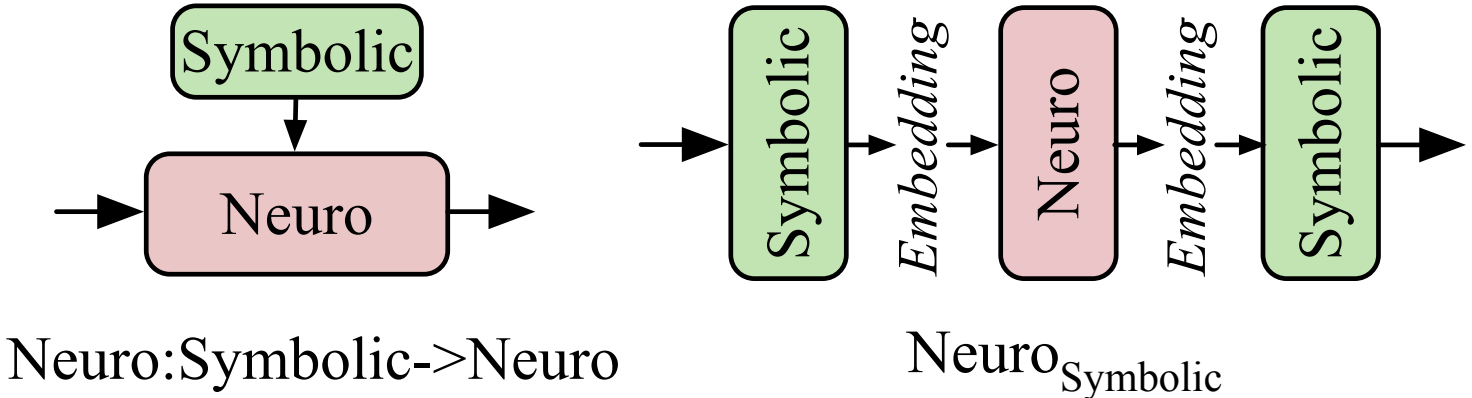
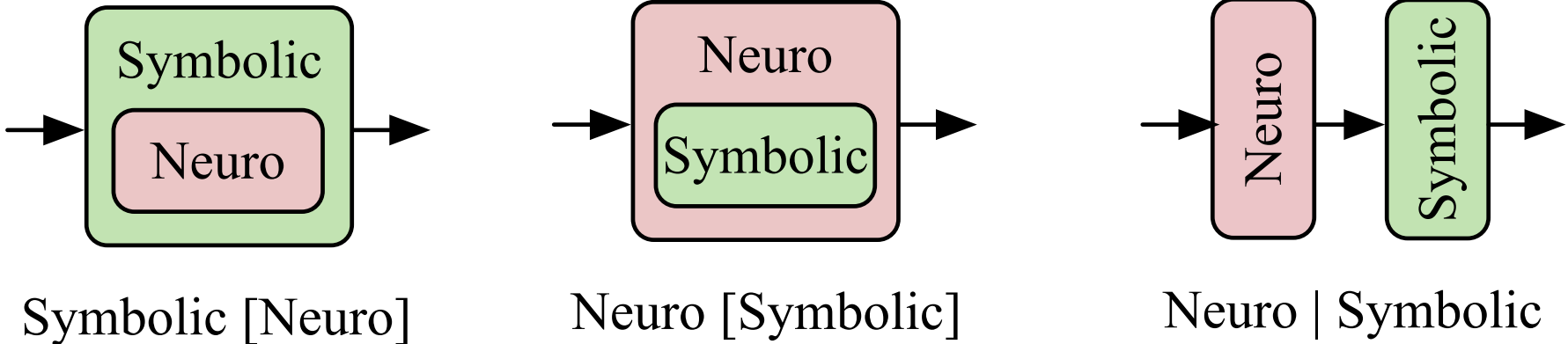


Neural Logical Machine

Neuro MLP, ConvNet, Transformer, etc

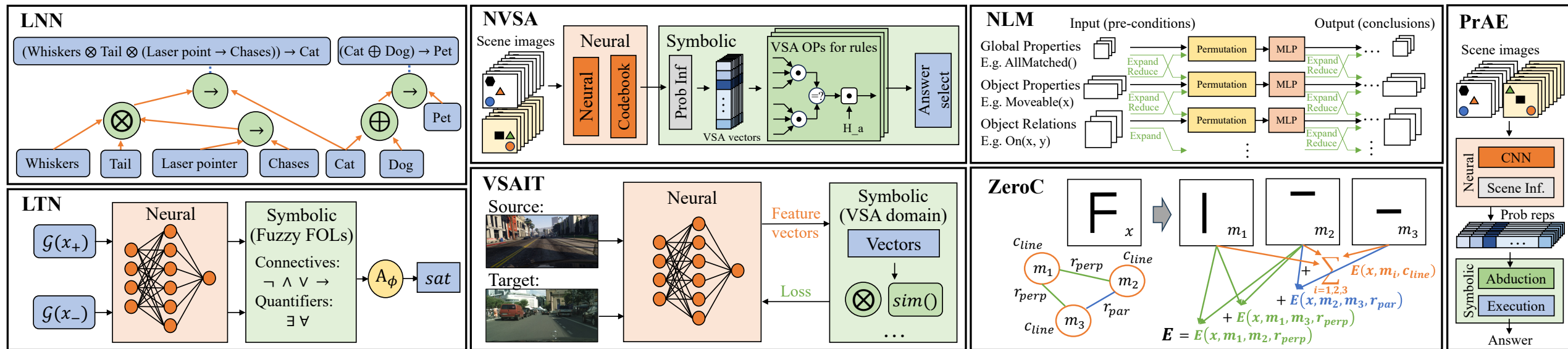
Symbolic Vector, Fuzzy logic, Knowledge graph, Decision tree, etc

Neuro-Symbolic AI Workload Category



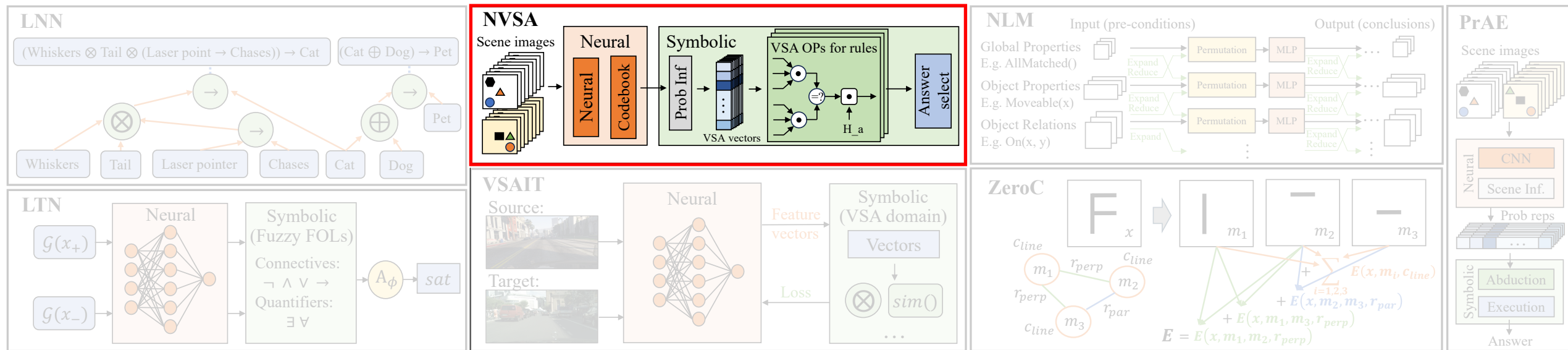
Inspired by Henry Kautz's terminology

Selected Neuro-Symbolic Workloads



Representative Neuro-Symbolic AI Workloads	Logic Neural Network [30]	Logic Tensor Network [34]	Neuro-Vector-Symbolic Architecture [4]	Vector Symbolic Architecture Image2Image Translation [7]	Neural Logic Machine [38]	Zero-shot Concept Recognition and Acquisition [37]	Probabilistic Abduction and Execution [23]	
Abbreviation	LNN	LTN	NVSA	VSAIT	NLM	ZeroC	PrAE	
Neuro-Symbolic Category	Neuro:Symbolic \rightarrow Neuro	Neuro _{Symbolic}	Neuro Symbolic	Neuro Symbolic	Neuro[Symbolic]	Neuro[Symbolic]	Neuro Symbolic	
Learning Approach	Supervised	Supervised/Unsupervised	Supervised/Unsupervised	Supervised	Supervised/Unsupervised	Supervised	Supervised/Unsupervised	
Deployment Scenario	Application	Learning and reasoning, Full theorem prover	Querying, learning, reasoning (relational and embedding learning, query answering)	Fluid intelligence, Abstract reasoning	Unpaired image-to-image translation	Relational reasoning, Decision making	Cross-domain classification and detection, Concept acquisition	Fluid intelligence, Spatial-temporal reasoning
	Advantage vs. Neural Model	Higher interoperability, resilience to incomplete knowledge, generalization	Higher data efficiency, comprehensibility, out-of-distribution generalization	Higher joint representations efficiency, abstract reasoning capability, transparency	Address semantic flipping and hallucinations issue in unpaired image translation tasks	Higher generalization, logic reasoning, deduction, explainability capability	Higher generalization, concept acquisition and recognition, compositionality capability	Higher generalization, transparency, interpretability, and robustness
	Dataset	LUBM benchmark [40], TPTP benchmark [41]	UCI [42], Leptograpsus crabs [43], DeepProbLog [44]	RAVEN [21], I-RAVEN [22], PGM [45]	GTA [47], Cityscapes [48], Google Maps dataset [49]	Family graph reasoning, sorting, path finding [46]	Abstraction reasoning [50], Hierarchical-concept corpus [51]	RAVEN [21], I-RAVEN [22], PGM [45]
Computation Pattern	Datatype	FP32	FP32	FP32	FP32	FP32	INT64	FP32
	Neuro	Graph	MLP	ConvNet	ConvNet	Sequential tensor	Energy-based network	ConvNet
	Symbolic	FOL/Logical operation	FOL/Logical operation	VSA/Vector operation	VSA/Vector operation	FOL/Logical operation	Graph, vector operation	VSA/Vector operation

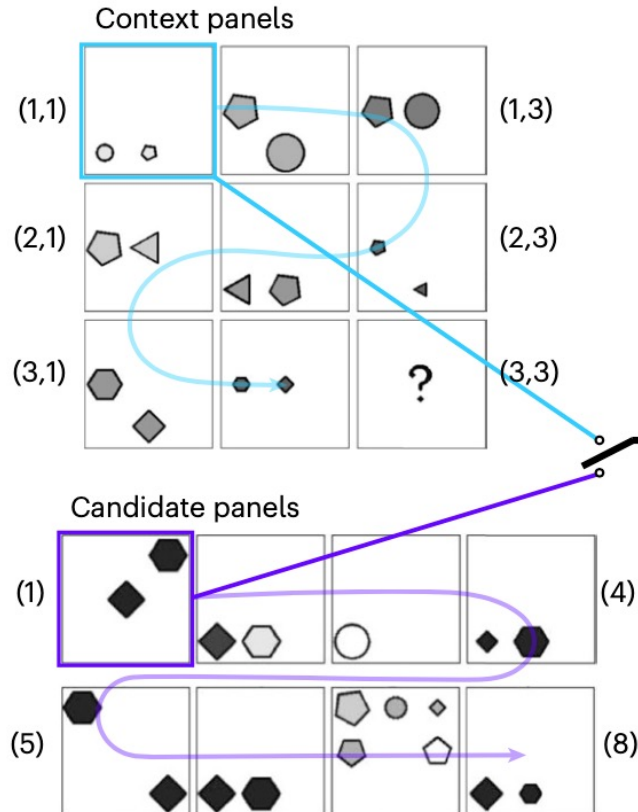
Selected Neuro-Symbolic Workloads



Representative Neuro-Symbolic AI Workloads	Logic Neural Network [30]	Logic Tensor Network [34]	Neuro-Vector-Symbolic Architecture [4]	Vector Symbolic Architecture Image2Image Translation [7]	Neural Logic Machine [38]	Zero-shot Concept Recognition and Acquisition [37]	Probabilistic Abduction and Execution [23]	
Abbreviation	LNN	LTN	NVSA	VSAIT	NLM	ZeroC	PrAE	
Neuro-Symbolic Category	Neuro:Symbolic \rightarrow Neuro	NeuroSymbolic	Neuro Symbolic	Neuro Symbolic	Neuro[Symbolic]	Neuro[Symbolic]	Neuro Symbolic	
Learning Approach	Supervised	Supervised/Unsupervised	Supervised/Unsupervised	Supervised	Supervised/Unsupervised	Supervised	Supervised/Unsupervised	
Deployment Scenario	Application	Learning and reasoning, Full theorem prover	Querying, learning, reasoning (relational and embedding learning, query answering)	Fluid intelligence, Abstract reasoning	Unpaired image-to-image translation	Relational reasoning, Decision making	Cross-domain classification and detection, Concept acquisition	Fluid intelligence, Spatial-temporal reasoning
	Advantage vs. Neural Model	Higher interoperability, resilience to incomplete knowledge, generalization	Higher data efficiency, comprehensibility, out-of-distribution generalization	Higher joint representations efficiency, abstract reasoning capability, transparency	Address semantic flipping and hallucinations issue in unpaired image translation tasks	Higher generalization, logic reasoning, deduction, explainability capability	Higher generalization, concept acquisition and recognition, compositionality capability	Higher generalization, transparency, interpretability, and robustness
	Dataset	LUBM benchmark [40], TPTP benchmark [41]	UCI [42], Leptograpsus crabs [43], DeepProbLog [44]	RAVEN [21], I-RAVEN [22], PGM [45]	GTA [47], Cityscapes [48], Google Maps dataset [49]	Family graph reasoning, sorting, path finding [46]	Abstraction reasoning [50], Hierarchical-concept corpus [51]	RAVEN [21], I-RAVEN [22], PGM [45]
Computation Pattern	Datatype	FP32	FP32	FP32	FP32	FP32	INT64	FP32
	Neuro	Graph	MLP	ConvNet	ConvNet	Sequential tensor	Energy-based network	ConvNet
	Symbolic	FOL/Logical operation	FOL/Logical operation	VSA/Vector operation	VSA/Vector operation	FOL/Logical operation	Graph, vector operation	VSA/Vector operation

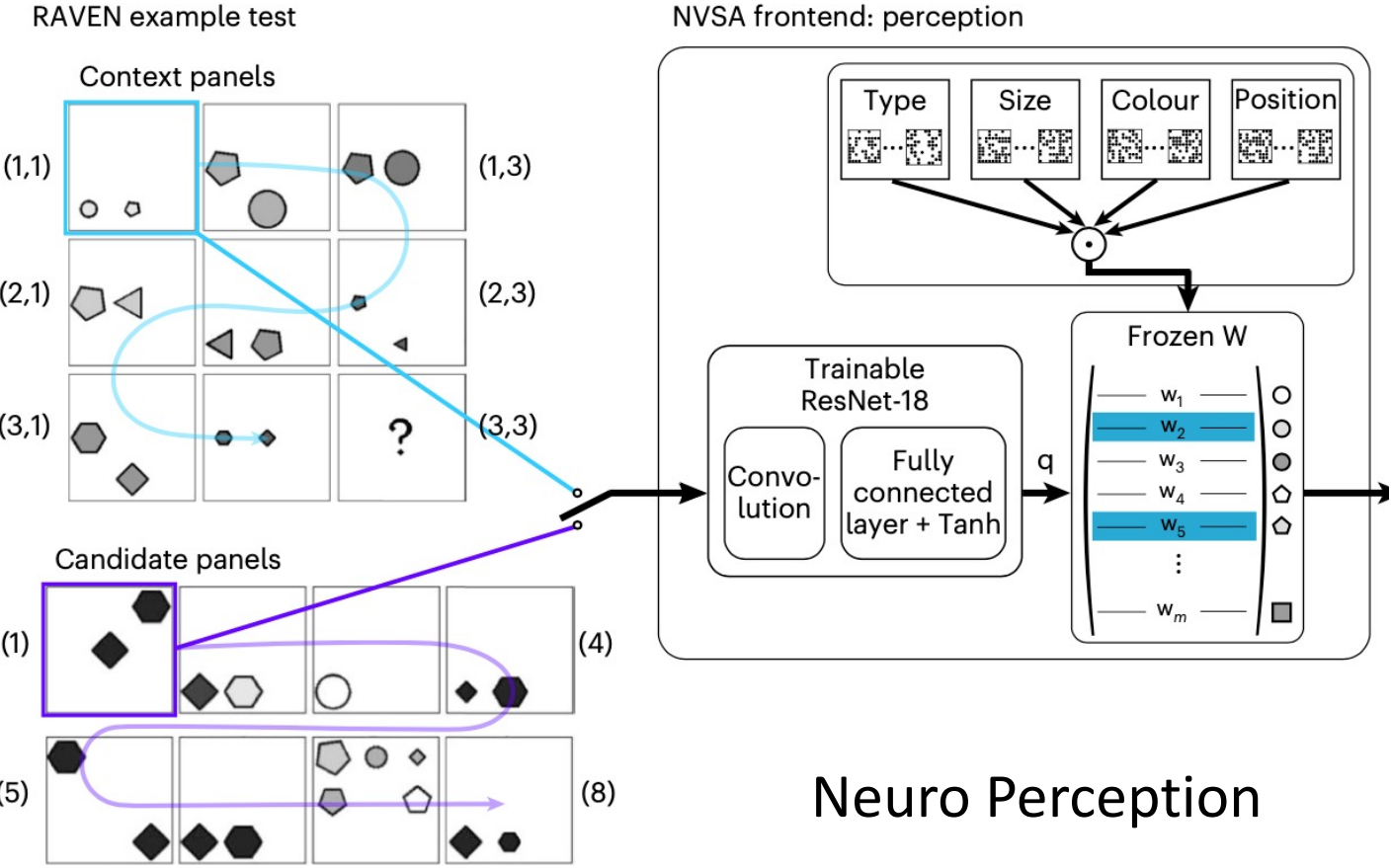
Example: Neuro-Vector-Symbolic Architecture

RAVEN example test



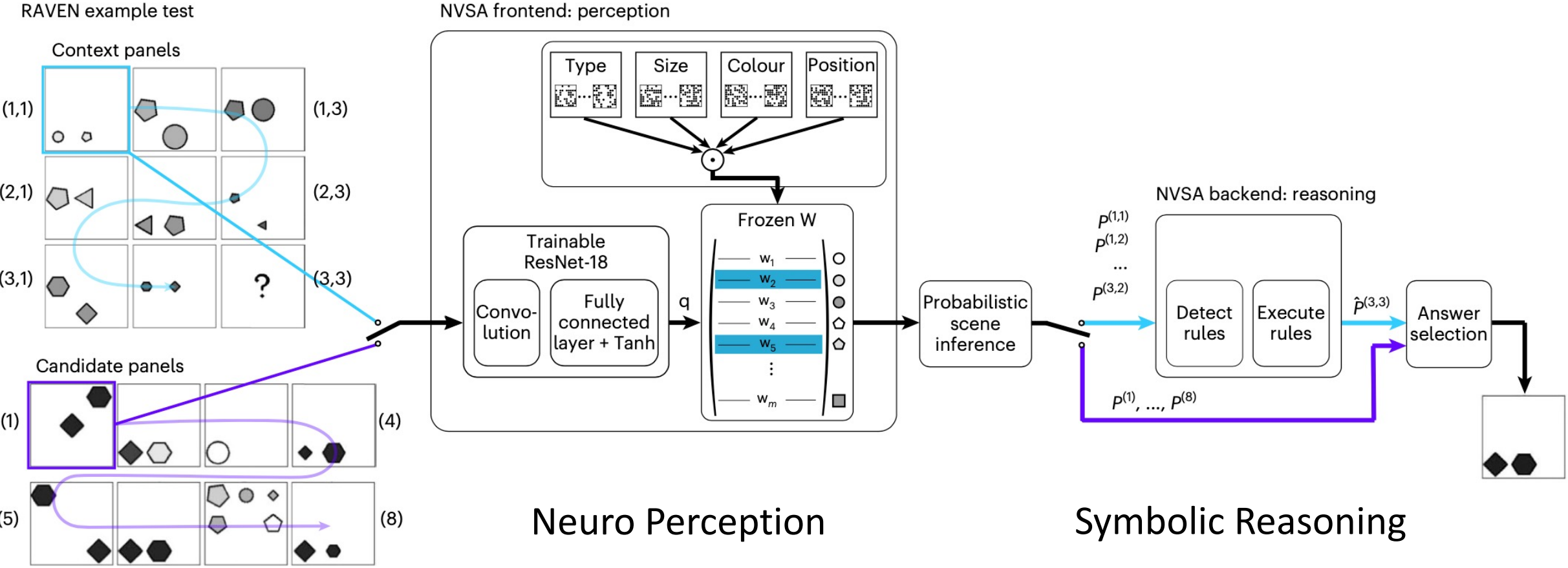
Hersche, et al. "A neuro-vector-symbolic architecture for solving Raven's progressive matrices". In Nature Machine Intelligence, 2023

Example: Neuro-Vector-Symbolic Architecture



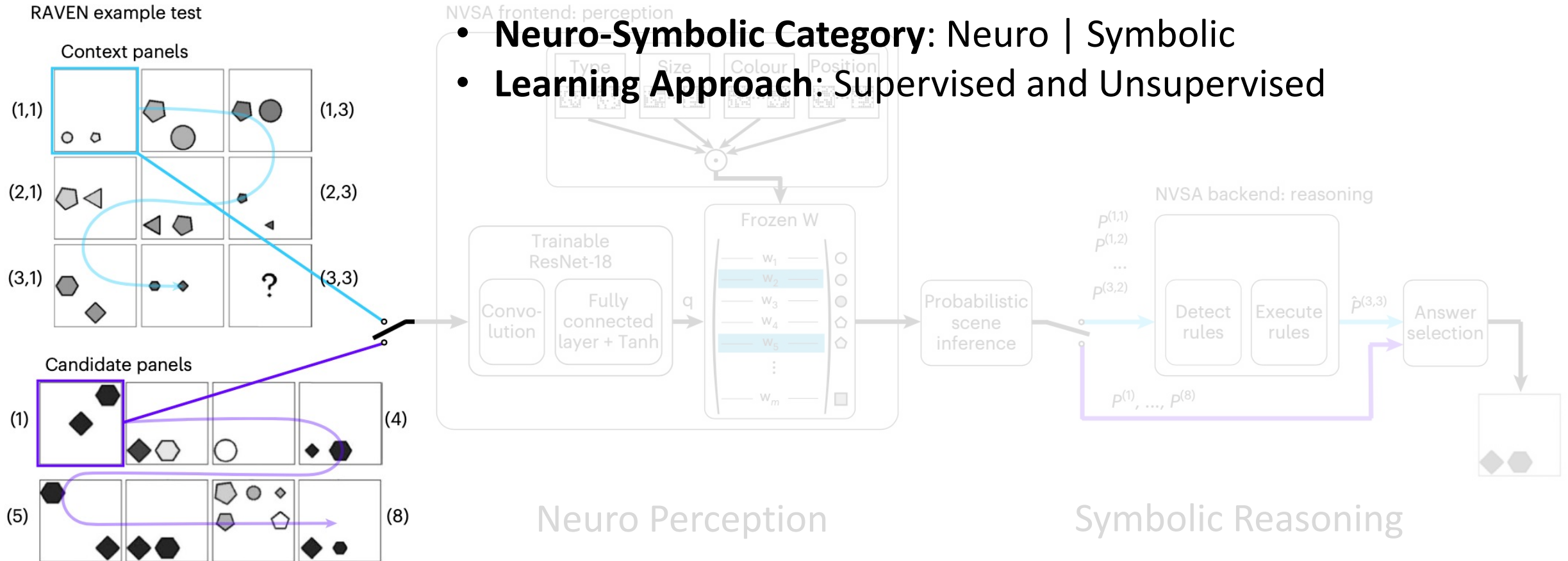
Hersche, et al. "A neuro-vector-symbolic architecture for solving Raven's progressive matrices". In Nature Machine Intelligence, 2023

Example: Neuro-Vector-Symbolic Architecture



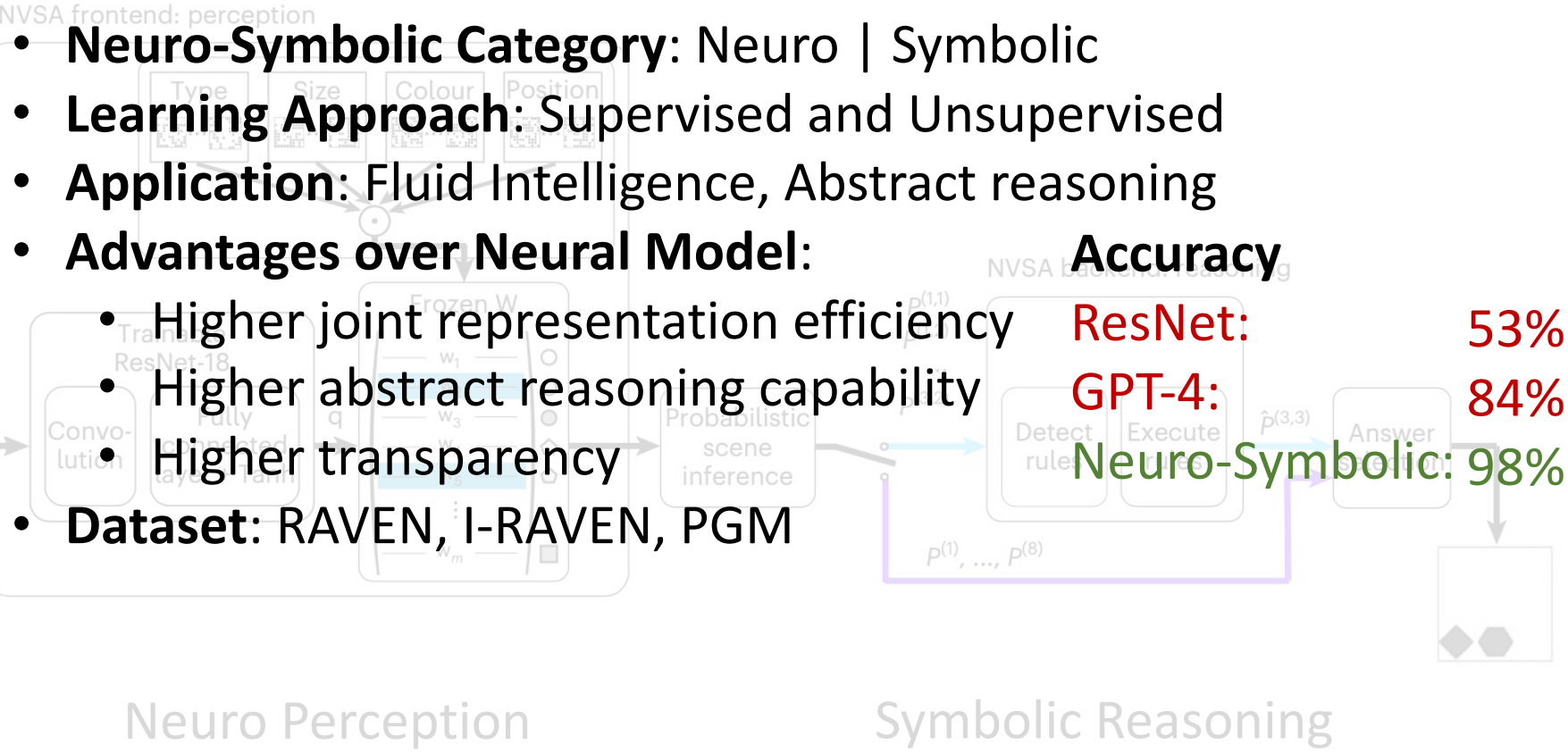
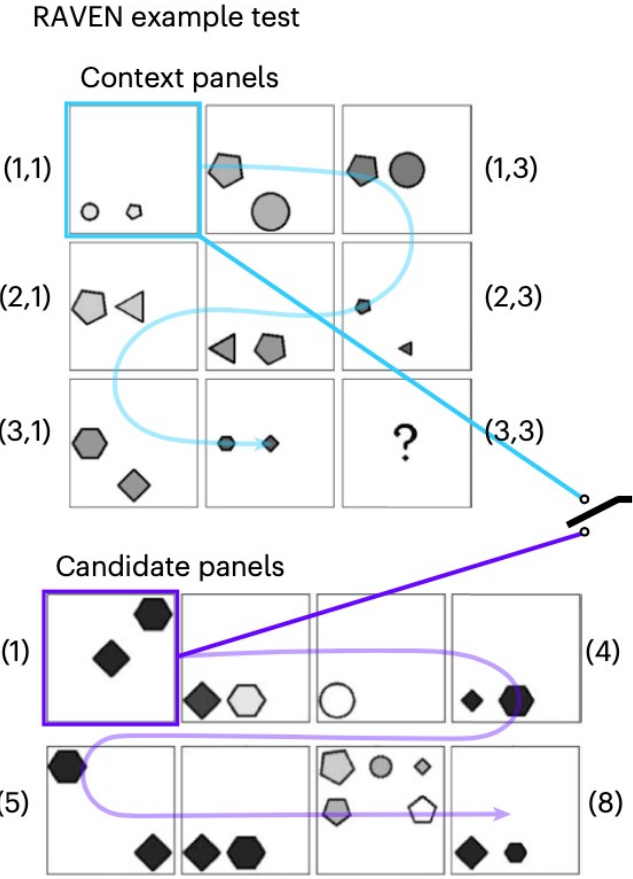
Hersche, et al. "A neuro-vector-symbolic architecture for solving Raven's progressive matrices". In Nature Machine Intelligence, 2023

Example: Neuro-Vector-Symbolic Architecture



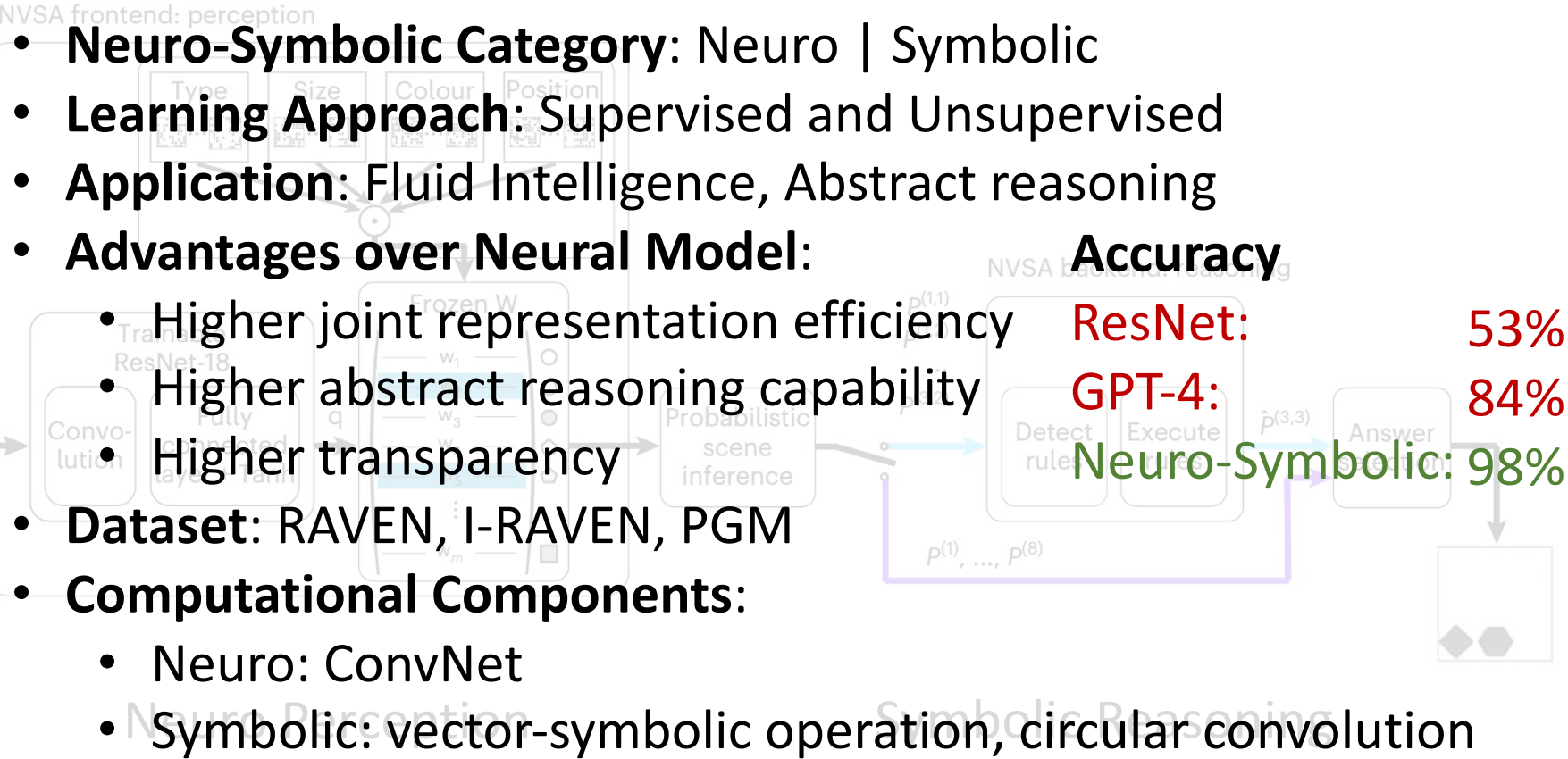
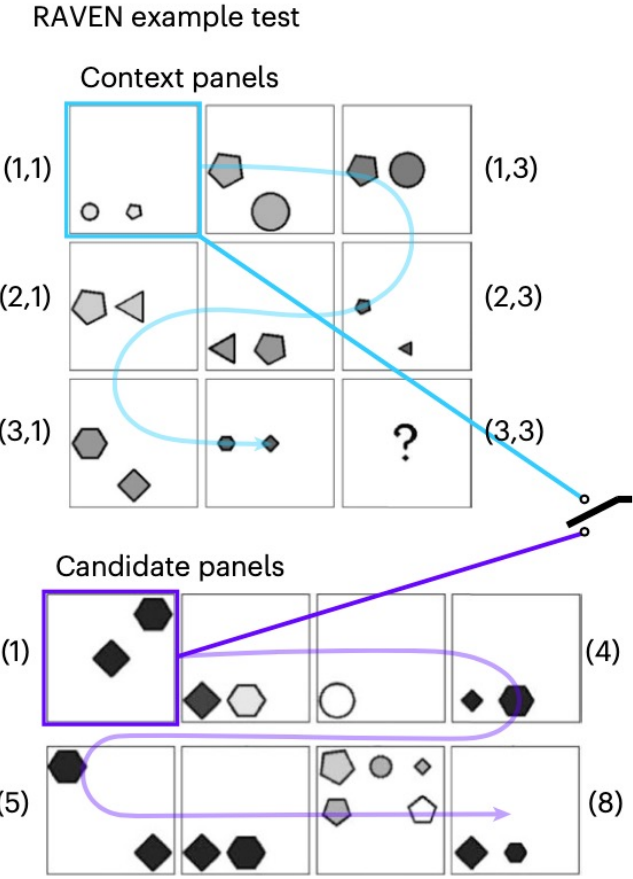
Hersche, et al. "A neuro-vector-symbolic architecture for solving Raven's progressive matrices". In Nature Machine Intelligence, 2023

Example: Neuro-Vector-Symbolic Architecture



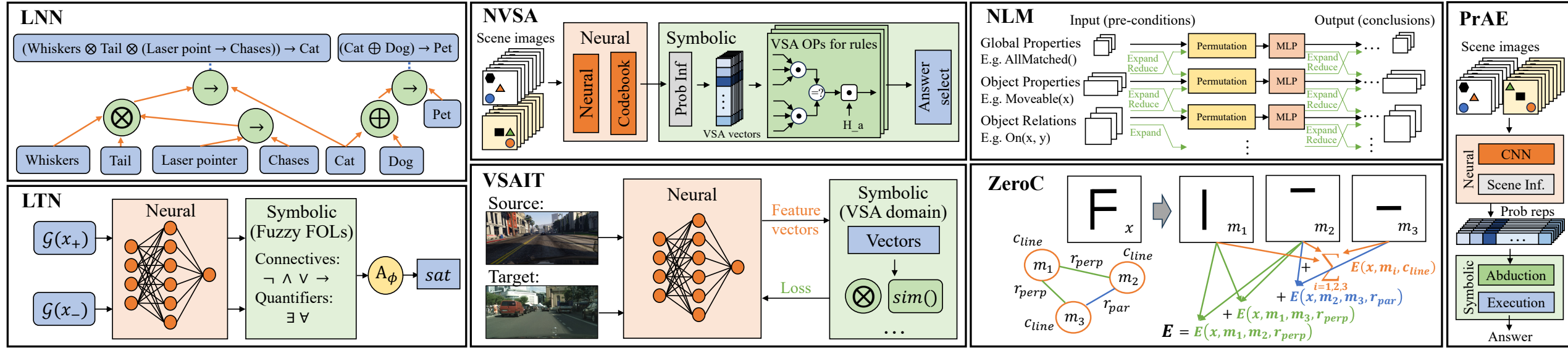
Hersche, et al. "A neuro-vector-symbolic architecture for solving Raven's progressive matrices". In Nature Machine Intelligence, 2023

Example: Neuro-Vector-Symbolic Architecture



Hersche, et al. "A neuro-vector-symbolic architecture for solving Raven's progressive matrices". In Nature Machine Intelligence, 2023

Selected Neuro-Symbolic Workloads



Representative Neuro-Symbolic AI Workloads	Logic Neural Network [30]	Logic Tensor Network [34]	Neuro-Vector-Symbolic Architecture [4]	Vector Symbolic Architecture Image2Image Translation [7]	Neural Logic Machine [38]	Zero-shot Concept Recognition and Acquisition [37]	Probabilistic Abduction and Execution [23]	
Abbreviation	LNN	LTN	NVSA	VSAIT	NLM	ZeroC	PrAE	
Neuro-Symbolic Category	Neuro:Symbolic \rightarrow Neuro	Neuro _{Symbolic}	Neuro Symbolic	Neuro Symbolic	Neuro[Symbolic]	Neuro[Symbolic]	Neuro Symbolic	
Learning Approach	Supervised	Supervised/Unsupervised	Supervised/Unsupervised	Supervised	Supervised/Unsupervised	Supervised	Supervised/Unsupervised	
Deployment Scenario	Application	Learning and reasoning, Full theorem prover	Querying, learning, reasoning (relational and embedding learning, query answering)	Fluid intelligence, Abstract reasoning	Unpaired image-to-image translation	Relational reasoning, Decision making	Cross-domain classification and detection, Concept acquisition	Fluid intelligence, Spatial-temporal reasoning
	Advantage vs. Neural Model	Higher interoperability, resilience to incomplete knowledge, generalization	Higher data efficiency, comprehensibility, out-of-distribution generalization	Higher joint representations efficiency, abstract reasoning capability, transparency	Address semantic flipping and hallucinations issue in unpaired image translation tasks	Higher generalization, logic reasoning, deduction, explainability capability	Higher generalization, concept acquisition and recognition, compositionality capability	Higher generalization, transparency, interpretability, and robustness
	Dataset	LUBM benchmark [40], TPTP benchmark [41]	UCI [42], Leptograpsus crabs [43], DeepProbLog [44]	RAVEN [21], I-RAVEN [22], PGM [45]	GTA [47], Cityscapes [48], Google Maps dataset [49]	Family graph reasoning, sorting, path finding [46]	Abstraction reasoning [50], Hierarchical-concept corpus [51]	RAVEN [21], I-RAVEN [22], PGM [45]
Computation Pattern	Datatype	FP32	FP32	FP32	FP32	FP32	INT64	FP32
	Neuro	Graph	MLP	ConvNet	ConvNet	Sequential tensor	Energy-based network	ConvNet
	Symbolic	FOL/Logical operation	FOL/Logical operation	VSA/Vector operation	VSA/Vector operation	FOL/Logical operation	Graph, vector operation	VSA/Vector operation

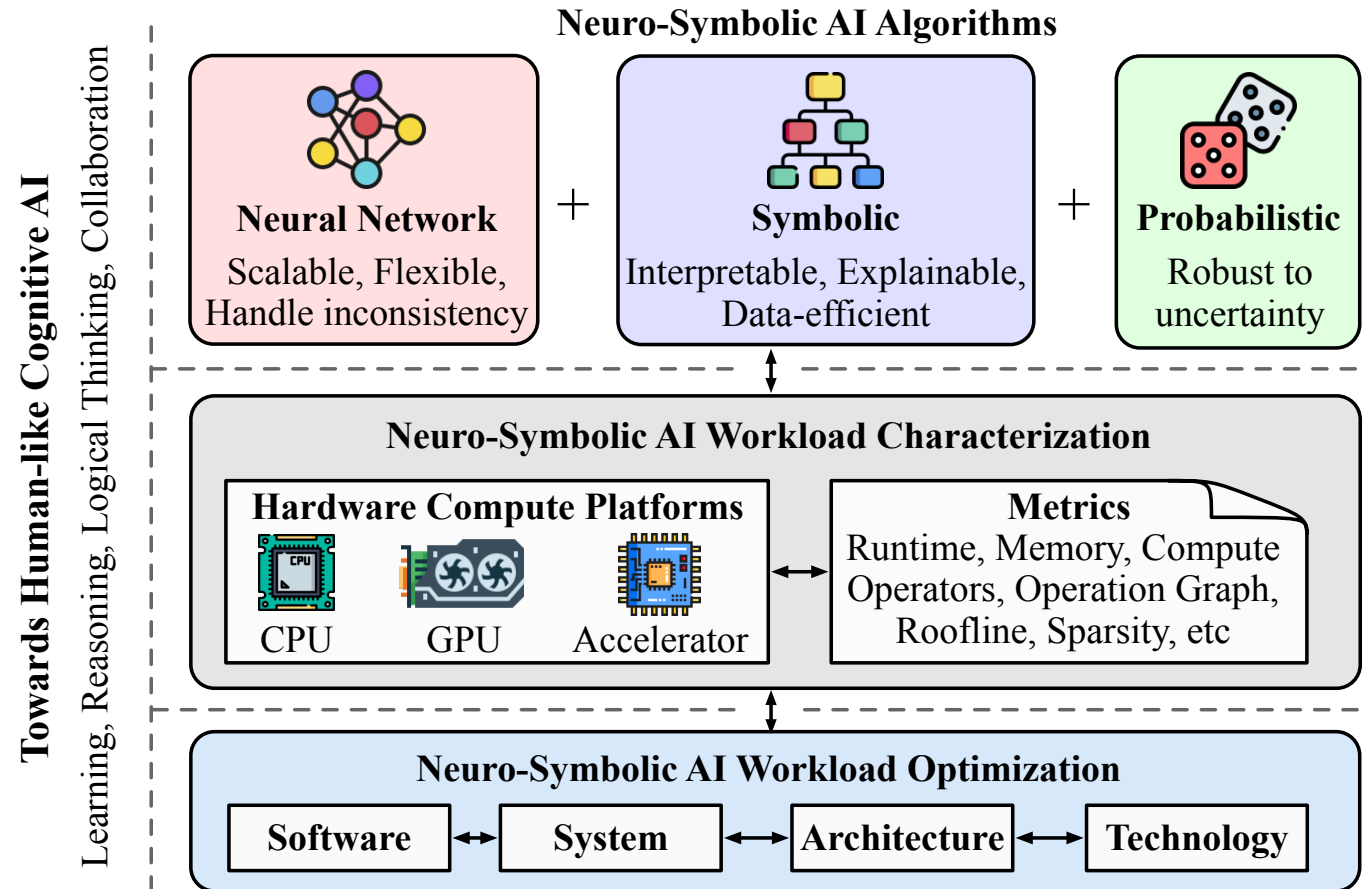
Objective of this Work

Workload and Characterization of Neuro-Symbolic AI

Categorize Neuro-Symbolic Algorithms

Understand Computational Behavior of Neuro-Symbolic Workloads

Identify Co-Design Opportunities



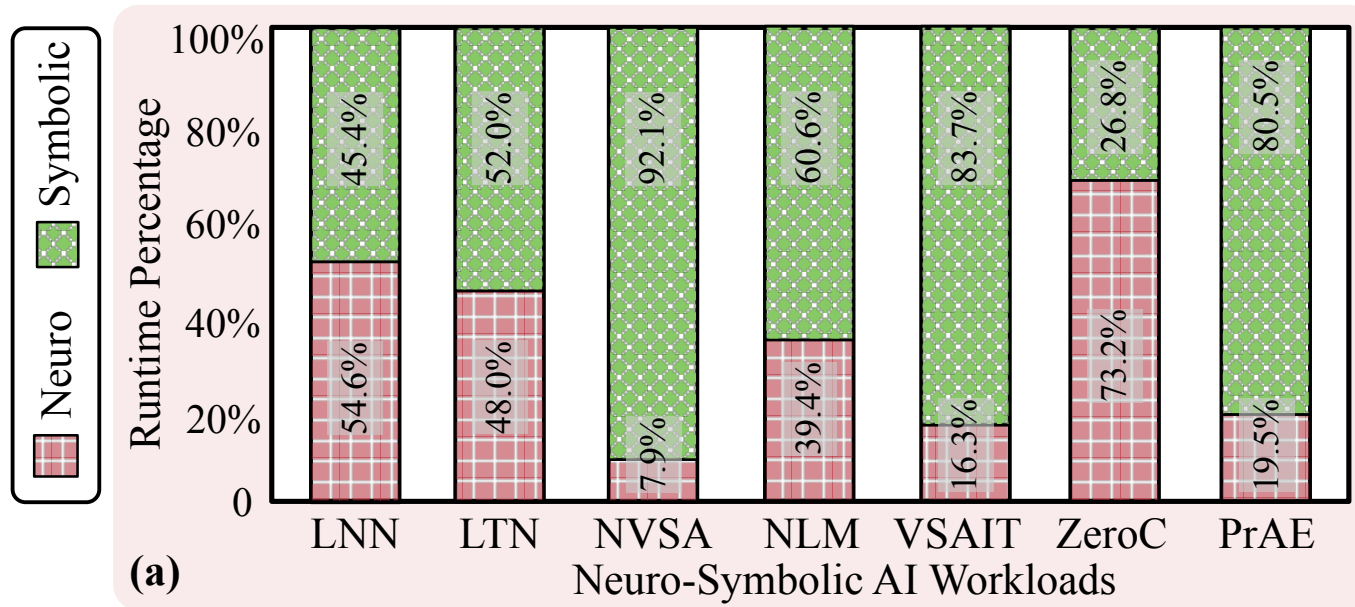
Neuro-Symbolic Workload Characterization

Profiling setup: CPU+GPU system, using pytorch profiler, seven neuro-symbolic workloads

Neuro-Symbolic Workload Characterization

Profiling setup: CPU+GPU system, using pytorch profiler, seven neuro-symbolic workloads

- End-to-end runtime latency analysis:

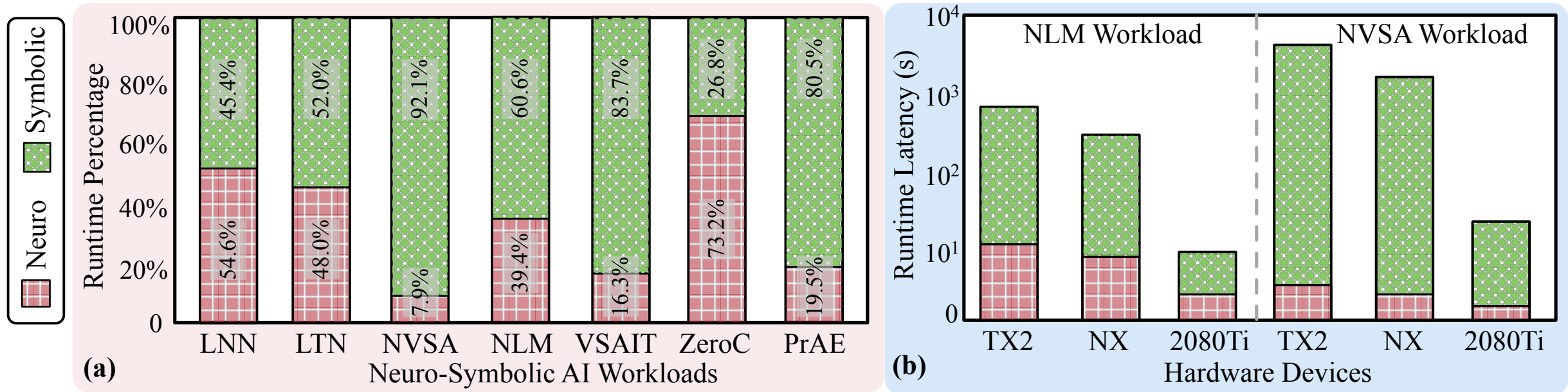


Neuro-symbolic workload exhibits high latency compared to neural models;
Symbolic component is processed inefficiently on off-the-shelf CPU/GPUs

Neuro-Symbolic Workload Characterization

Profiling setup: CPU+GPU system, using pytorch profiler, seven neuro-symbolic workloads

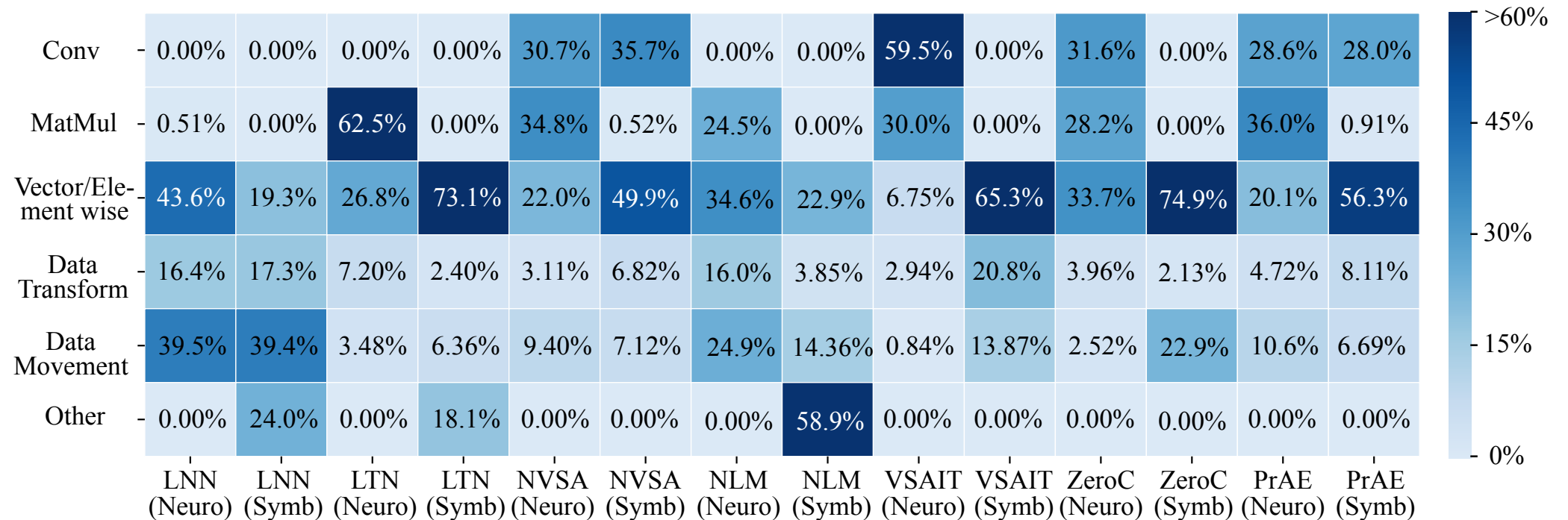
- End-to-end runtime latency analysis:



Neuro-symbolic workload exhibits high latency compared to neural models;
Symbolic component is processed inefficiently on off-the-shelf CPU/GPUs

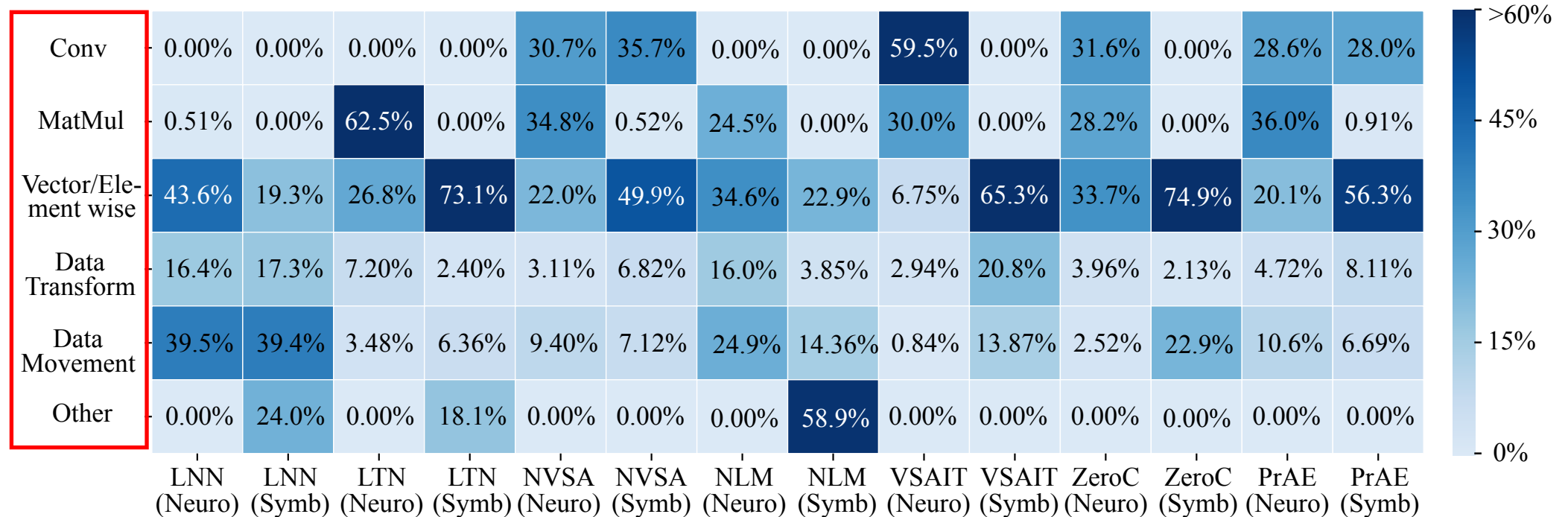
Neuro-Symbolic Workload Characterization

- Compute operator analysis:



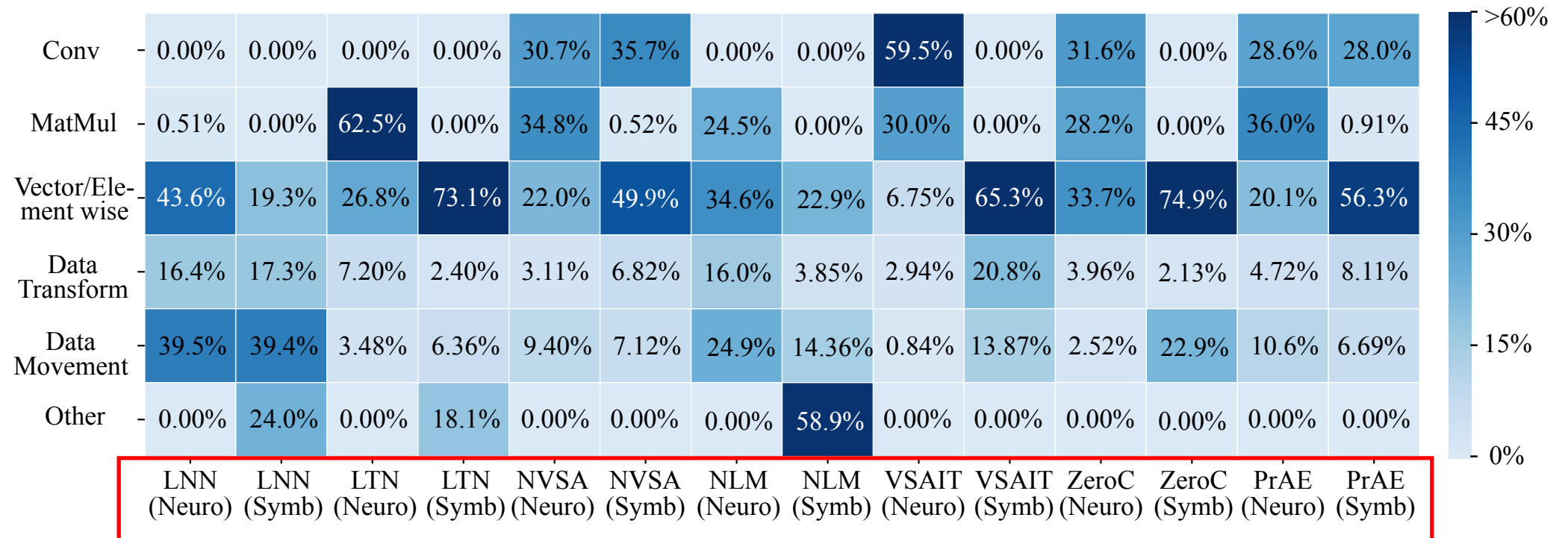
Neuro-Symbolic Workload Characterization

- Compute operator analysis:



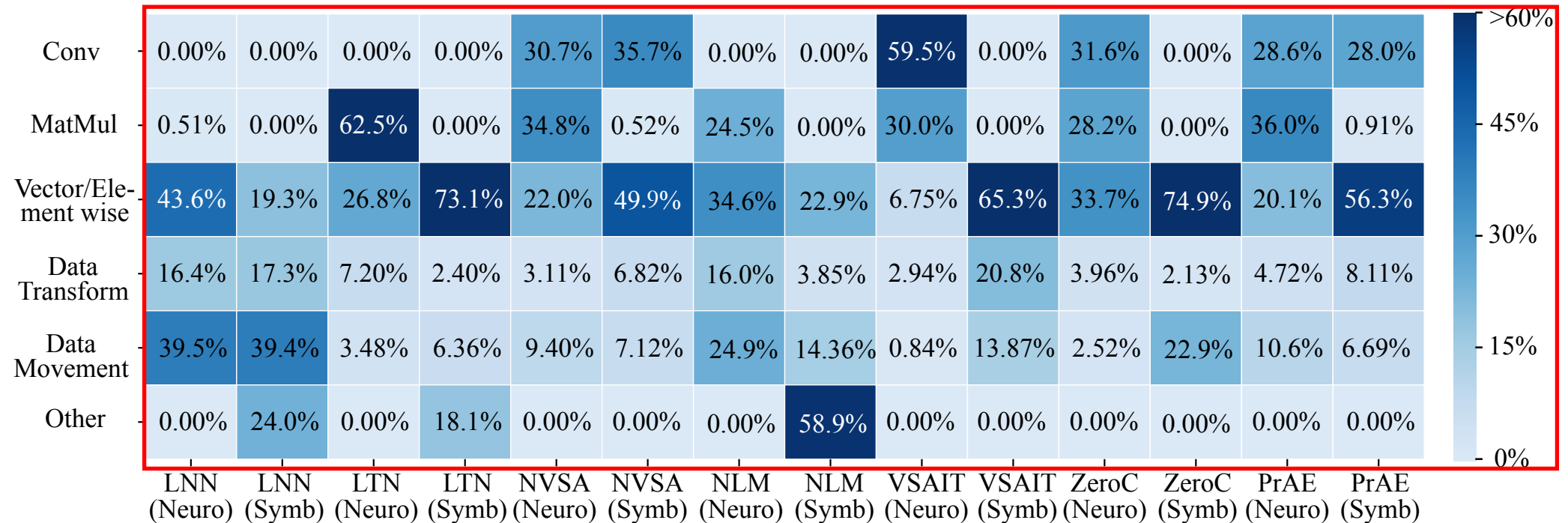
Neuro-Symbolic Workload Characterization

- Compute operator analysis:



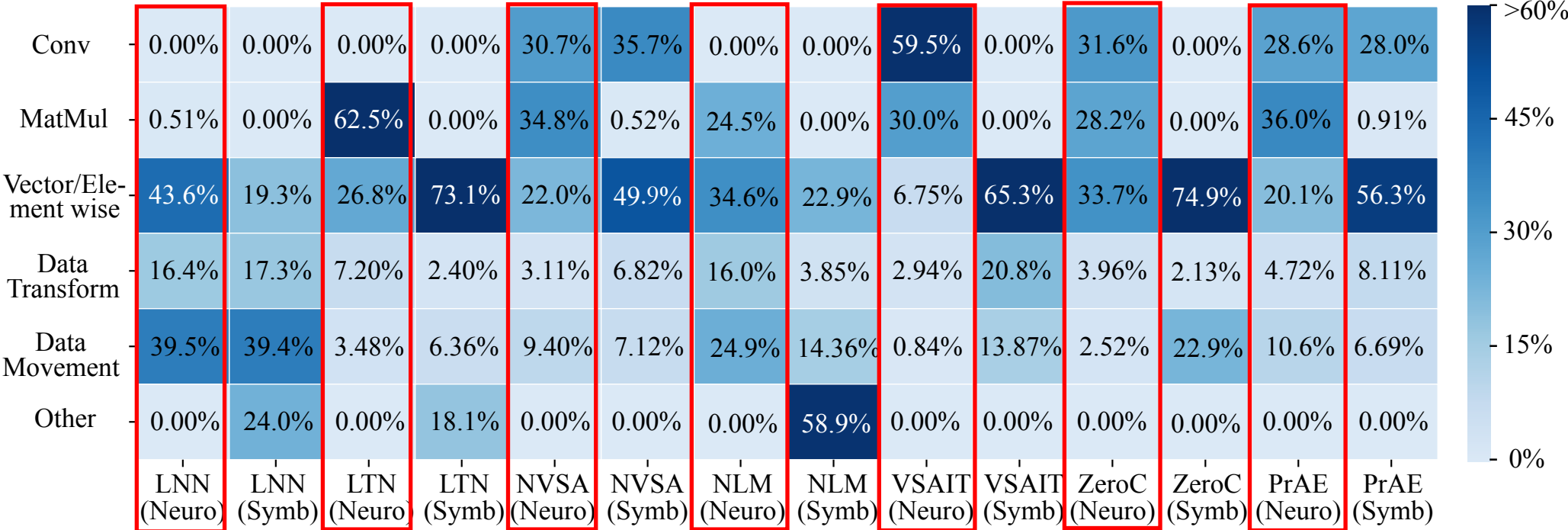
Neuro-Symbolic Workload Characterization

- Compute operator analysis:



Neuro-Symbolic Workload Characterization

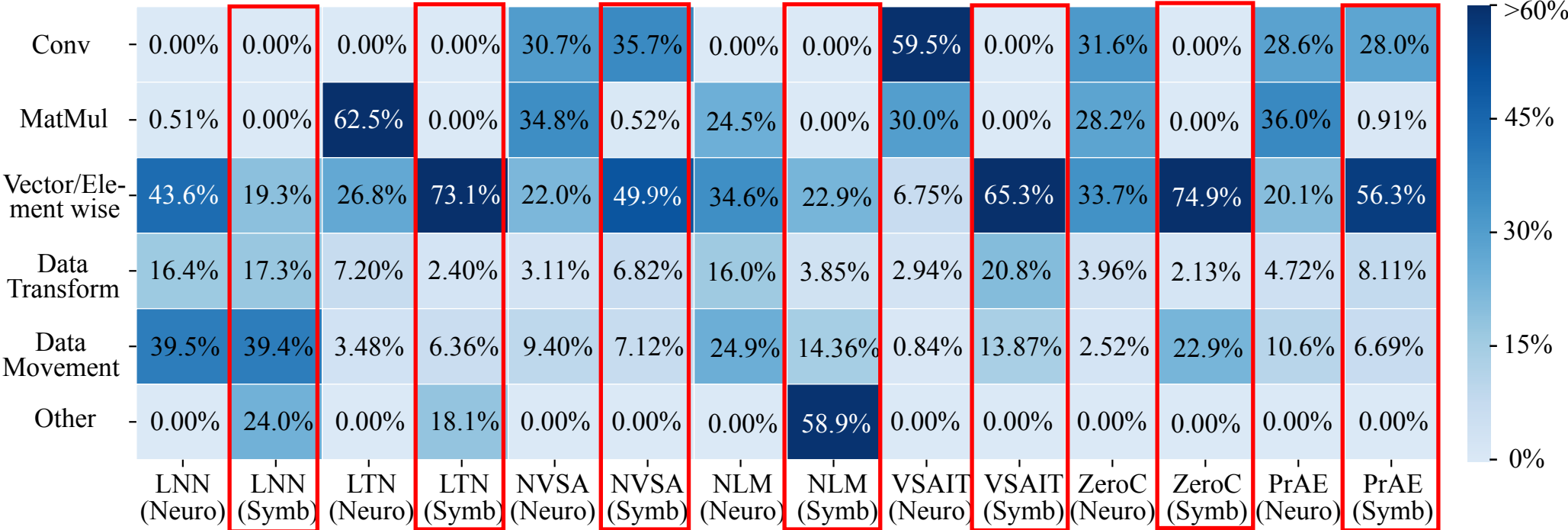
- Compute operator analysis:



Neural dominated by MatMul and Conv;

Neuro-Symbolic Workload Characterization

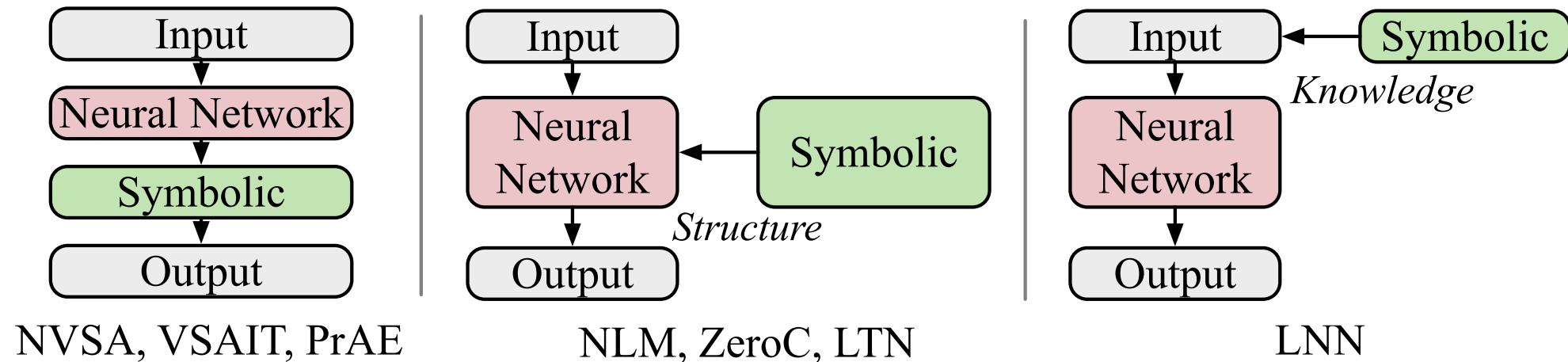
- Compute operator analysis:



Neural dominated by MatMul and Conv; Symbolic dominated by vector/element/logical operations;

Neuro-Symbolic Workload Characterization

- Data Dependence Graph analysis:



Neural dominated by MatMul and Conv; Symbolic dominated by vector/element/logical operations; Complex control flow of neuro-symbolic interaction

Neuro-Symbolic Workload Characterization

	Neuro Kernel		Symbolic Kernel	
	segmm_nn	relu_nn	vectorized	elementwise
Runtime Percentage (%)				
Compute Throughput (%)				
ALU Utilization (%)				
L1 Cache Hit Rate (%)				
L2 Cache Hit Rate (%)				
L1 Cache Throughput (%)				
L2 Cache Throughput (%)				
DRAM BW Utilization (%)				

Why system Inefficiency?

Neuro-Symbolic Workload Characterization

	Neuro Kernel		Symbolic Kernel	
	segmm_nn	relu_nn	vectorized	elementwise
Runtime Percentage (%)	18.2	10.4	37.5	12.4
Compute Throughput (%)	95.1	92.9	3.0	2.3
ALU Utilization (%)	90.1	48.3	5.9	4.5
L1 Cache Hit Rate (%)				
L2 Cache Hit Rate (%)				
L1 Cache Throughput (%)				
L2 Cache Throughput (%)				
DRAM BW Utilization (%)				

Symbolic exhibits low ALU utilization,

Neuro-Symbolic Workload Characterization

	Neuro Kernel		Symbolic Kernel	
	segmm_nn	relu_nn	vectorized	elementwise
Runtime Percentage (%)	18.2	10.4	37.5	12.4
Compute Throughput (%)	95.1	92.9	3.0	2.3
ALU Utilization (%)	90.1	48.3	5.9	4.5
L1 Cache Hit Rate (%)	1.6	51.6	29.5	33.3
L2 Cache Hit Rate (%)	86.8	65.5	48.6	34.3
L1 Cache Throughput (%)				
L2 Cache Throughput (%)				
DRAM BW Utilization (%)				

Symbolic exhibits low ALU utilization, low cache hit rate,

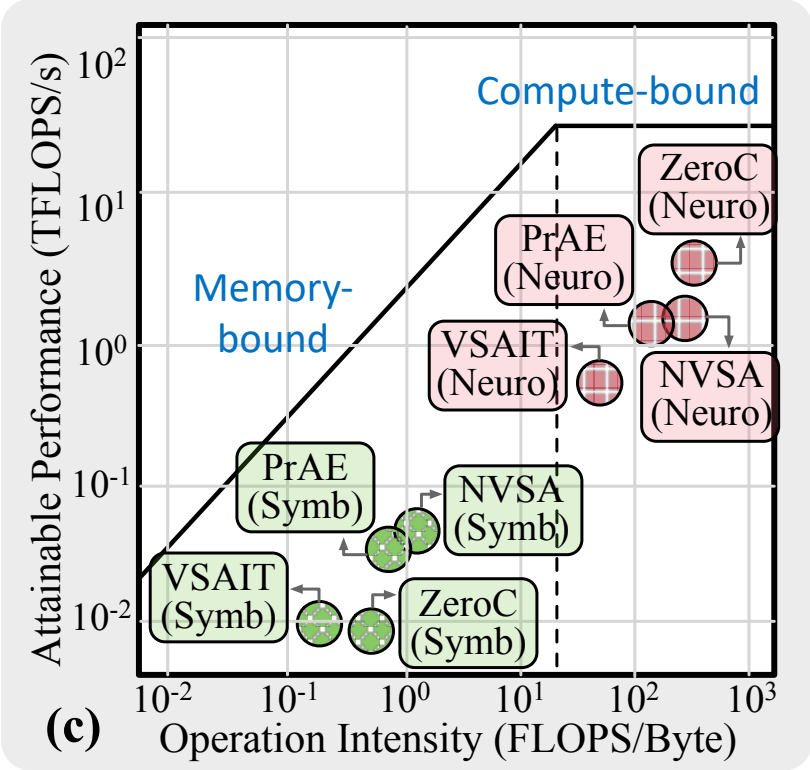
Neuro-Symbolic Workload Characterization

	Neuro Kernel		Symbolic Kernel	
	segmm_nn	relu_nn	vectorized	elementwise
Runtime Percentage (%)	18.2	10.4	37.5	12.4
Compute Throughput (%)	95.1	92.9	3.0	2.3
ALU Utilization (%)	90.1	48.3	5.9	4.5
L1 Cache Hit Rate (%)	1.6	51.6	29.5	33.3
L2 Cache Hit Rate (%)	86.8	65.5	48.6	34.3
L1 Cache Throughput (%)	79.7	82.6	28.4	10.8
L2 Cache Throughput (%)	19.2	17.5	29.8	22.8
DRAM BW Utilization (%)	14.9	24.2	90.9	78.4

Symbolic exhibits low ALU utilization, low cache hit rate, massive data transfer, resulting in hardware underutilization and inefficiency

Neuro-Symbolic Workload Characterization

	Neuro Kernel		Symbolic Kernel	
	segmm_nn	relu_nn	vectorized	elementwise
Runtime Percentage (%)	18.2	10.4	37.5	12.4
Compute Throughput (%)	95.1	92.9	3.0	2.3
ALU Utilization (%)	90.1	48.3	5.9	4.5
L1 Cache Hit Rate (%)	1.6	51.6	29.5	33.3
L2 Cache Hit Rate (%)	86.8	65.5	48.6	34.3
L1 Cache Throughput (%)	79.7	82.6	28.4	10.8
L2 Cache Throughput (%)	19.2	17.5	29.8	22.8
DRAM BW Utilization (%)	14.9	24.2	90.9	78.4



Neuro operations are compute-bounded, symbolic operations are memory-bounded.

Neural Network vs. Neuro-Symbolic

	Neural Network	Neuro-Symbolic

Neural Network vs. Neuro-Symbolic

	Neural Network	Neuro-Symbolic
Runtime	[Neural Network] < [Neural-Symbolic]	

Neural Network vs. Neuro-Symbolic

	Neural Network	Neuro-Symbolic
Runtime	[Neural Network] < [Neural-Symbolic]	
Compute Kernels	Neural kernels (Conv, MatMul, etc)	Heterogenous neural and symbolic kernels (vector, element, MatMul, graph, logic, etc)

Neural Network vs. Neuro-Symbolic

	Neural Network	Neuro-Symbolic
Runtime	[Neural Network] < [Neural-Symbolic]	
Compute Kernels	Neural kernels (Conv, MatMul, etc)	Heterogenous neural and symbolic kernels (vector, element, MatMul, graph, logic, etc)
Hardware Efficiency	Efficient on GPU/TPU	Inefficient on CPU/GPU/TPU (low ALU utilization, low L1 cache hit rate, high data movement, etc)

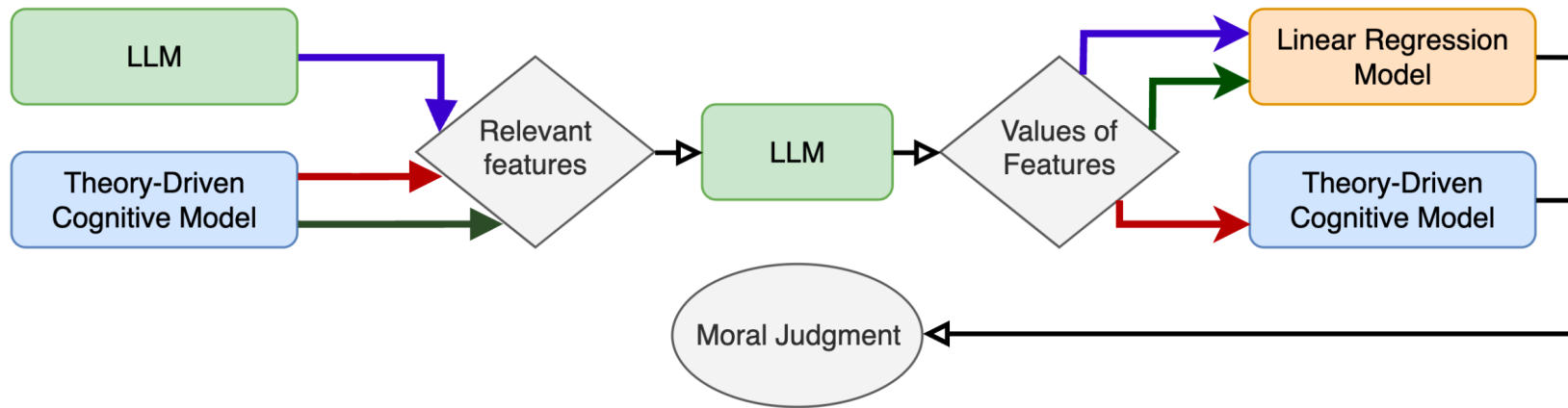
Neural Network vs. Neuro-Symbolic

	Neural Network	Neuro-Symbolic
Runtime	[Neural Network] < [Neural-Symbolic]	
Compute Kernels	Neural kernels (Conv, MatMul, etc)	Heterogenous neural and symbolic kernels (vector, element, MatMul, graph, logic, etc)
Hardware Efficiency	Efficient on GPU/TPU	Inefficient on CPU/GPU/TPU (low ALU utilization, low L1 cache hit rate, high data movement, etc)
System Bound	Compute-bound / Memory-bound	Memory-bound

Neural Network vs. Neuro-Symbolic

	Neural Network	Neuro-Symbolic
Runtime	[Neural Network] < [Neural-Symbolic]	
Compute Kernels	Neural kernels (Conv, MatMul, etc)	Heterogenous neural and symbolic kernels (vector, element, MatMul, graph, logic, etc)
Hardware Efficiency	Efficient on GPU/TPU	Inefficient on CPU/GPU/TPU (low ALU utilization, low L1 cache hit rate, high data movement, etc)
System Bound	Compute-bound / Memory-bound	Memory-bound
Dataflow	Simple flow control, High parallelism	Complex flow control, Low parallelism

Looking Ahead: LLM + Neurosymbolic



Scenario:
Imagine that a stranger will give Hank one thousand dollars to break all the windows in his neighbor's house without his neighbor's permission. Hank carries out the stranger's request.

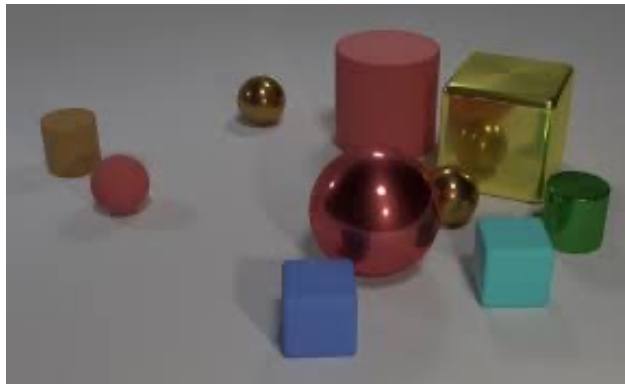
Towards safe and trustworthy AI System:
LLM + cognitive symbolic model for human moral judgment

Looking Ahead: Challenge and Opportunity

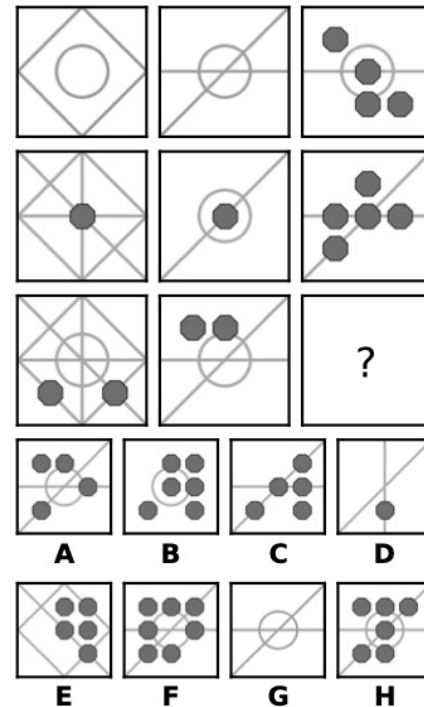
Data



Lack of cognitive datasets



CLEVRER Dataset



RAVEN Dataset

Looking Ahead: Challenge and Opportunity

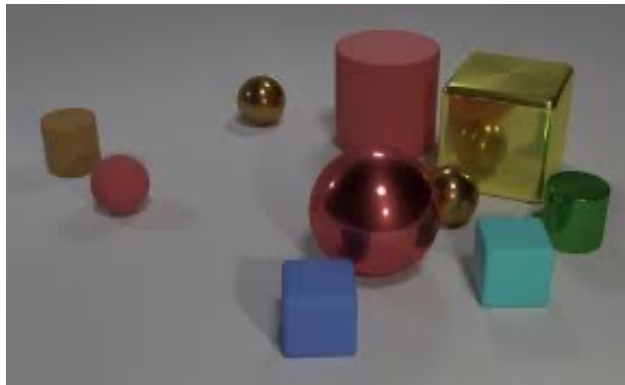
Data



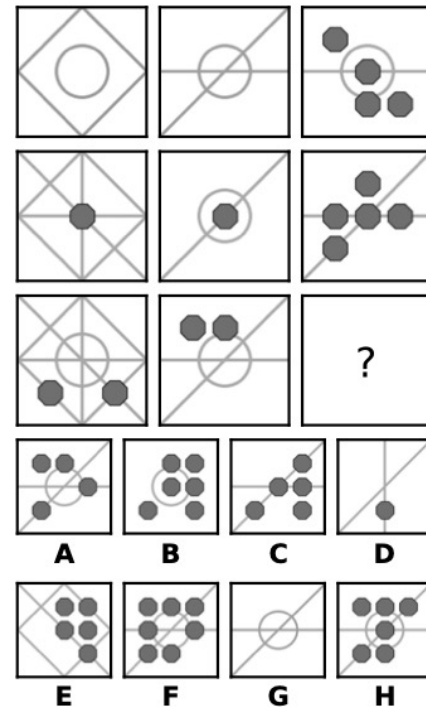
Lack of cognitive datasets



Building ImageNet-like NSAI datasets



CLEVRER Dataset



RAVEN Dataset



Human-like AI

Metacognition
Interpretability
Deductive Reasoning
Systematicity
Compositionality
Counterfactual thinking

...

Looking Ahead: Challenge and Opportunity

Data



Lack of cognitive datasets



Building ImageNet-like NSAI datasets

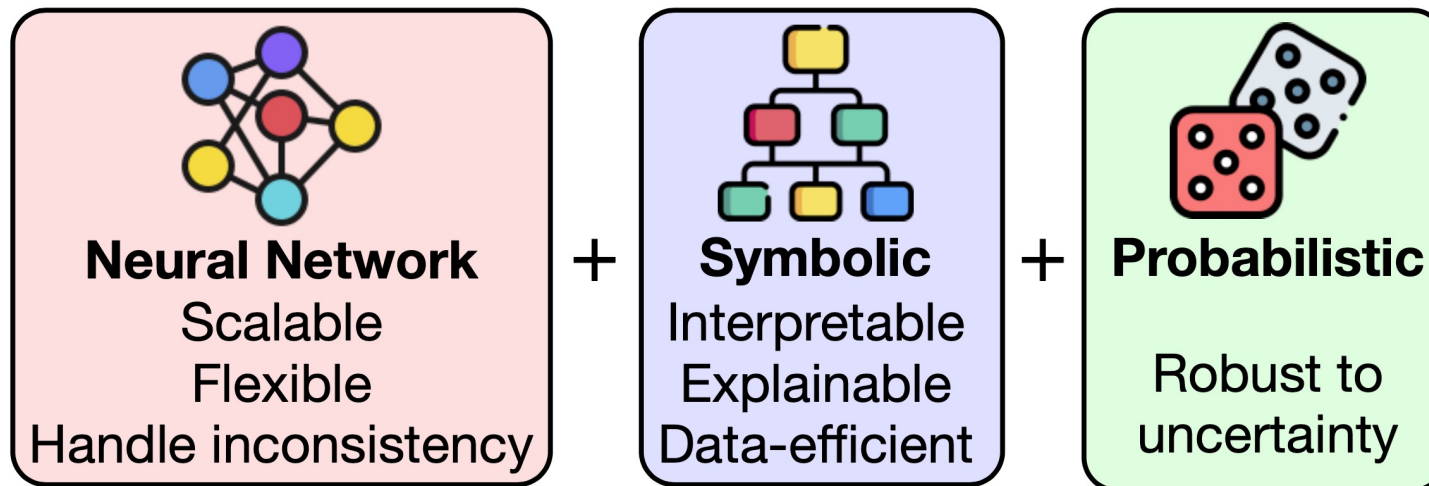
Model



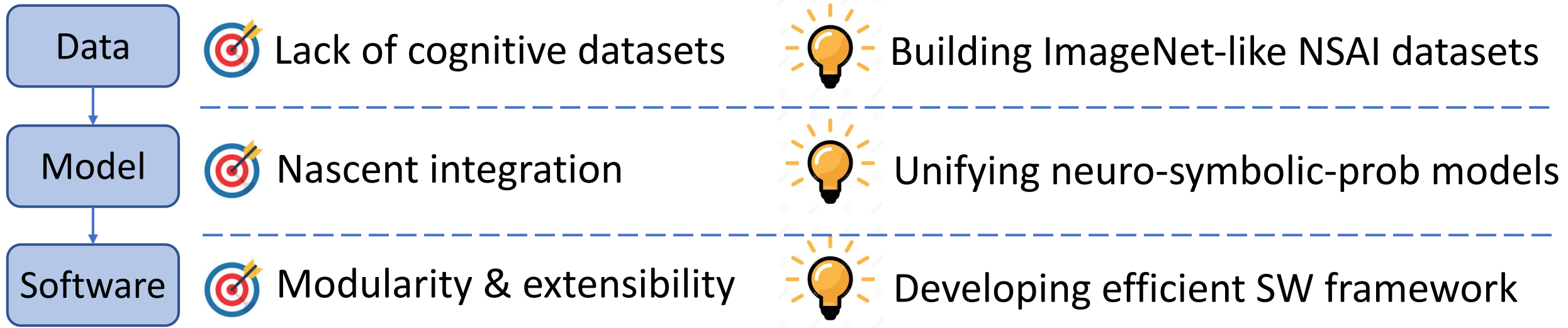
Nascent integration



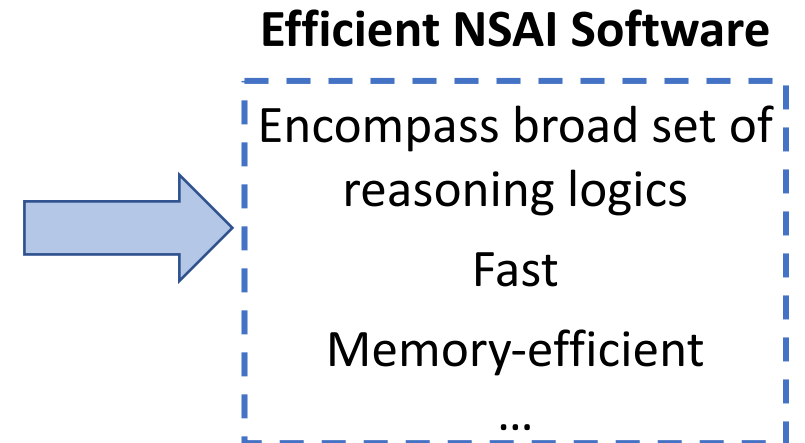
Unifying neuro-symbolic-prob models



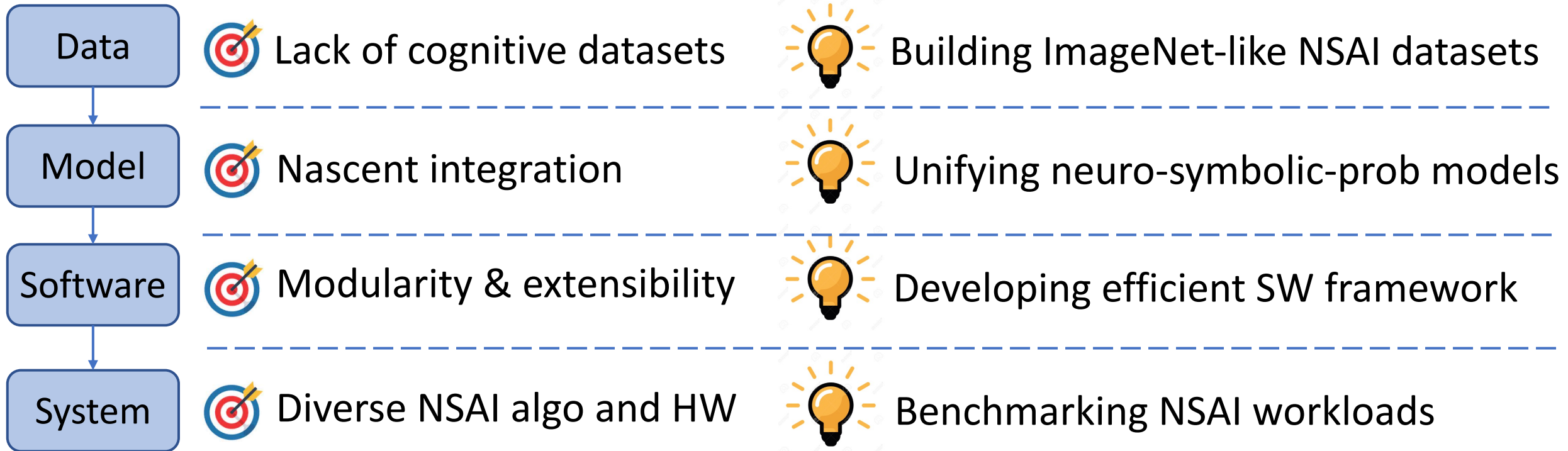
Looking Ahead: Challenge and Opportunity



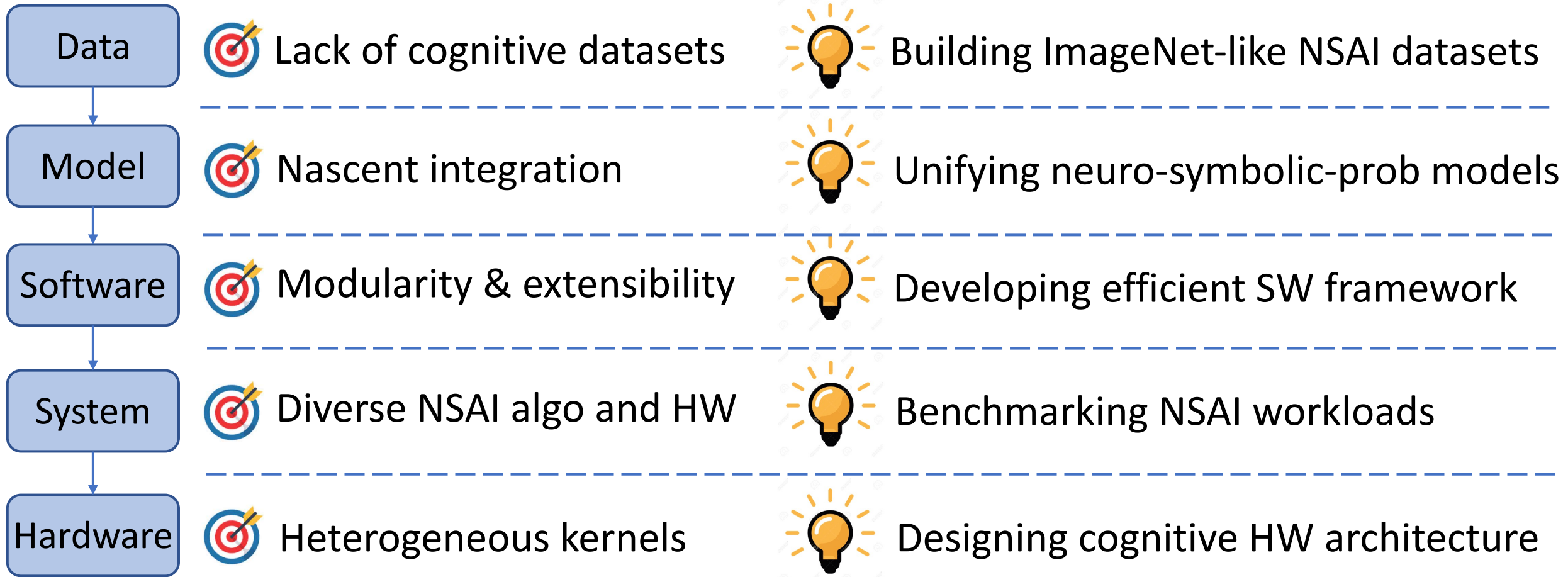
Underlying Operations	Examples
Fuzzy logic (LTN)	$F = \forall x(isCarnivor(s)) \rightarrow (isMammal(x))$ $\{isCarnivor(s):[0, 1], isMammal(x):[1, 0]\} \rightarrow F = [1, 0]$
Mul and Add (NVSA)	$X_i \in \{+1, -1\}^d \rightarrow (X_i \cdot X_j) / (X_i + X_j)$
Pre-defined objects (NSVQA)	<code>equal_color: (entry, entry) → Boolean</code> <code>equal_integer: (number, number) → Boolean</code>



Looking Ahead: Challenge and Opportunity



Looking Ahead: Challenge and Opportunity



Summary

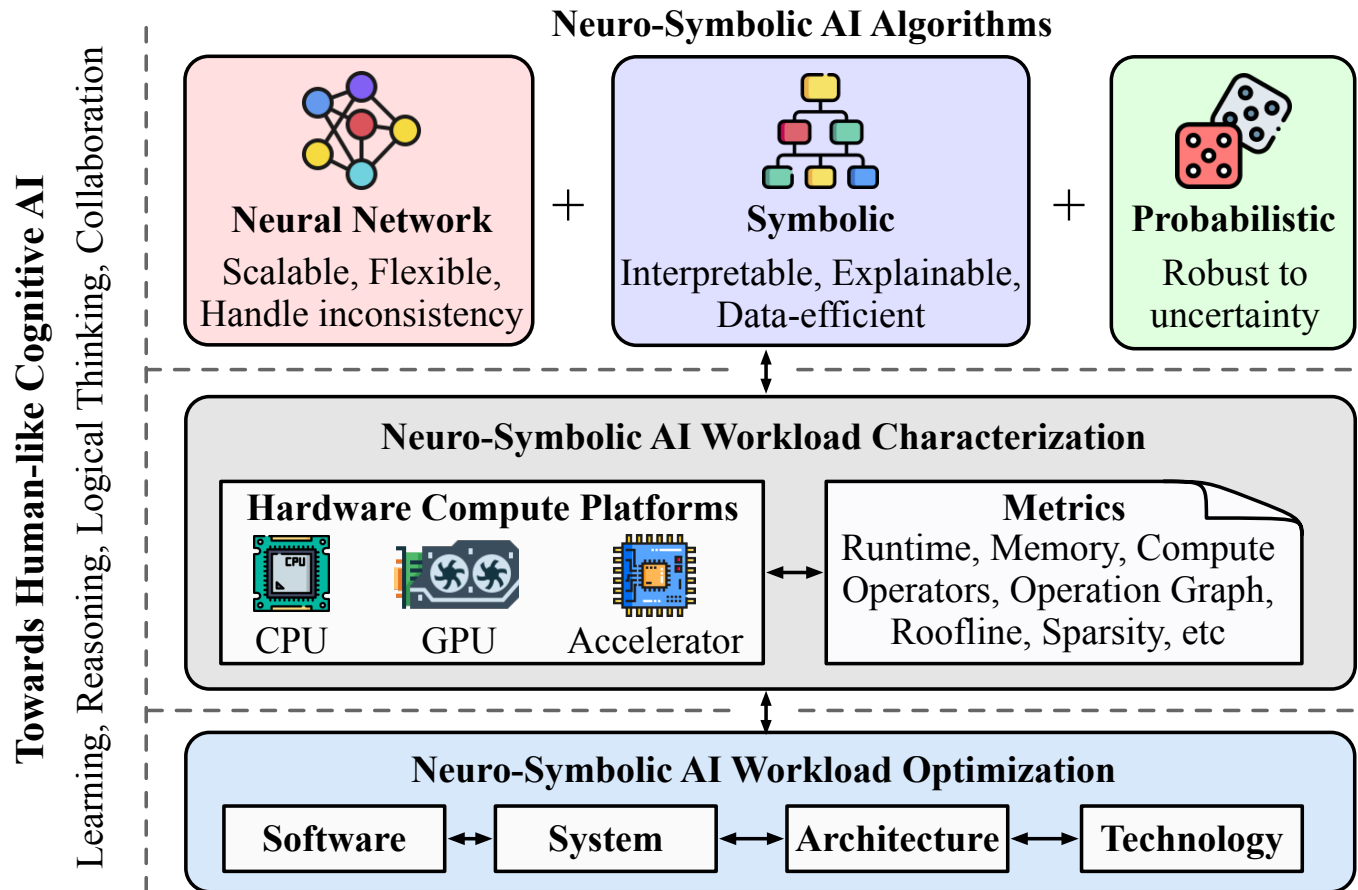
Thank you!!

Workload and Characterization of Neuro-Symbolic AI

Categorize Neuro-Symbolic Algorithms

Understand Computational Behavior of Neuro-Symbolic Workloads

Identify Co-Design Opportunities



Towards Cognitive AI Systems: Workload and Characterization of Neuro-Symbolic AI

Zishen Wan¹, Che-Kai Liu¹, Hanchen Yang¹, Ritik Raj¹, Chaojian Li¹, Haoran You¹, Yonggan Fu¹, Cheng Wan¹, Ananda Samajdar², Yingyan (Celine) Lin¹, Tushar Krishna¹, Arijit Raychowdhury¹

¹ *Georgia Institute of Technology, GA* ² *IBM Research, NY*

Email: zishenwan@gatech.edu

Web: <https://zishenwan.github.io>