

# Reasoning with Small & Large Models

Vaishak Belle, *University of Edinburgh*

DL excels at (scalable)  
pattern recognition but ...

*structured reasoning, causality vs.  
correlation, data hungry, explainability,  
abstraction, ...*

# Neuro-symbolic AI: response to the limitations of DL

Pragmatically, formalisms and frameworks that *combine, enhance* or *support* neural networks with reasoning and symbolic emergence mechanisms

**"Scaling is all you need" hypothesis vs the  
need for (hybrid) neurosymbolic systems**

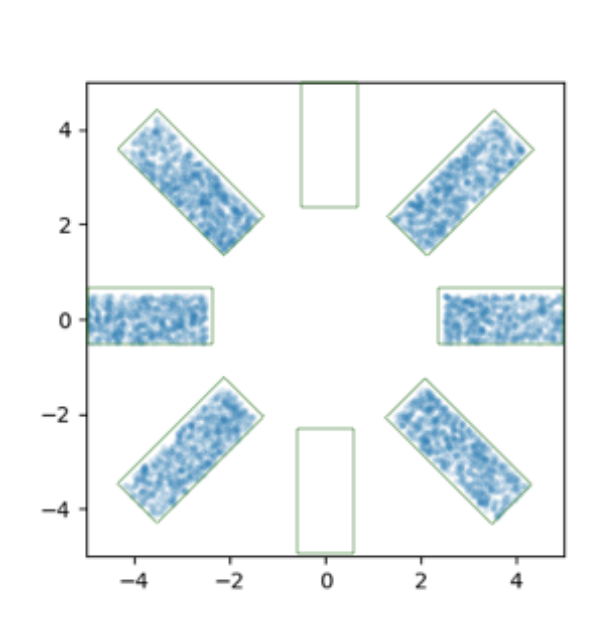
# Examples

Knowledge graphs (vs RAGs)

**Constrained prediction** (including RL)

Program execution & learning

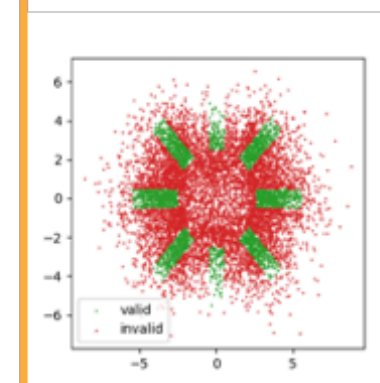
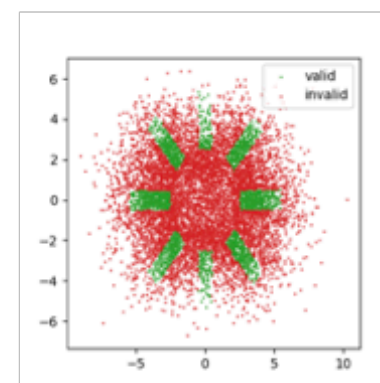
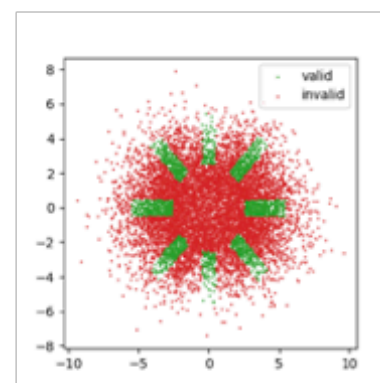
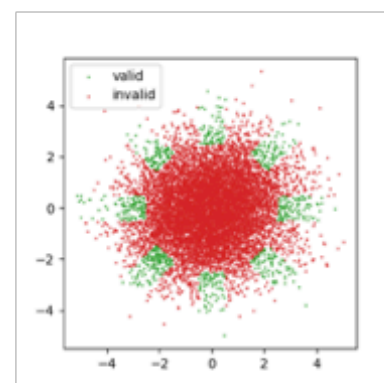
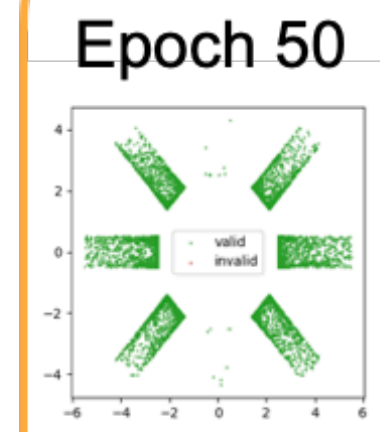
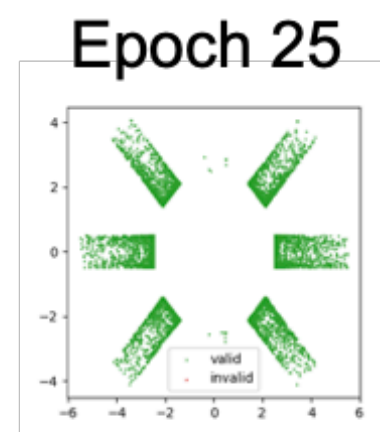
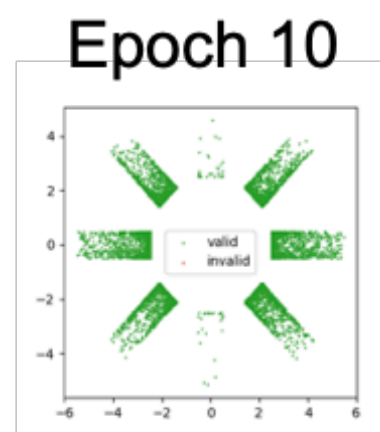
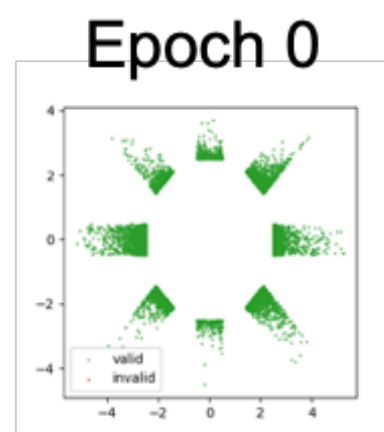
**LLM + executor**



MultiplexNet

Baseline VAE

Invalid samples high!



# LLMs as Encoders

## Premise

There are four persons. Everyone is visible to others. Each person draws a card, face unrevealed (red or black). Cara's card is shown to Vasiliki. Cara's card is shown to Conrad. Jennifer's card is shown to Conrad. Vasiliki's card is shown to Cara. It is publicly announced that someone picked a red card. It is publicly announced that Vasiliki knows whether someone picked a red card.

.....

## Hypothesis

Cara can now know whether Conrad picked a red card.

.....

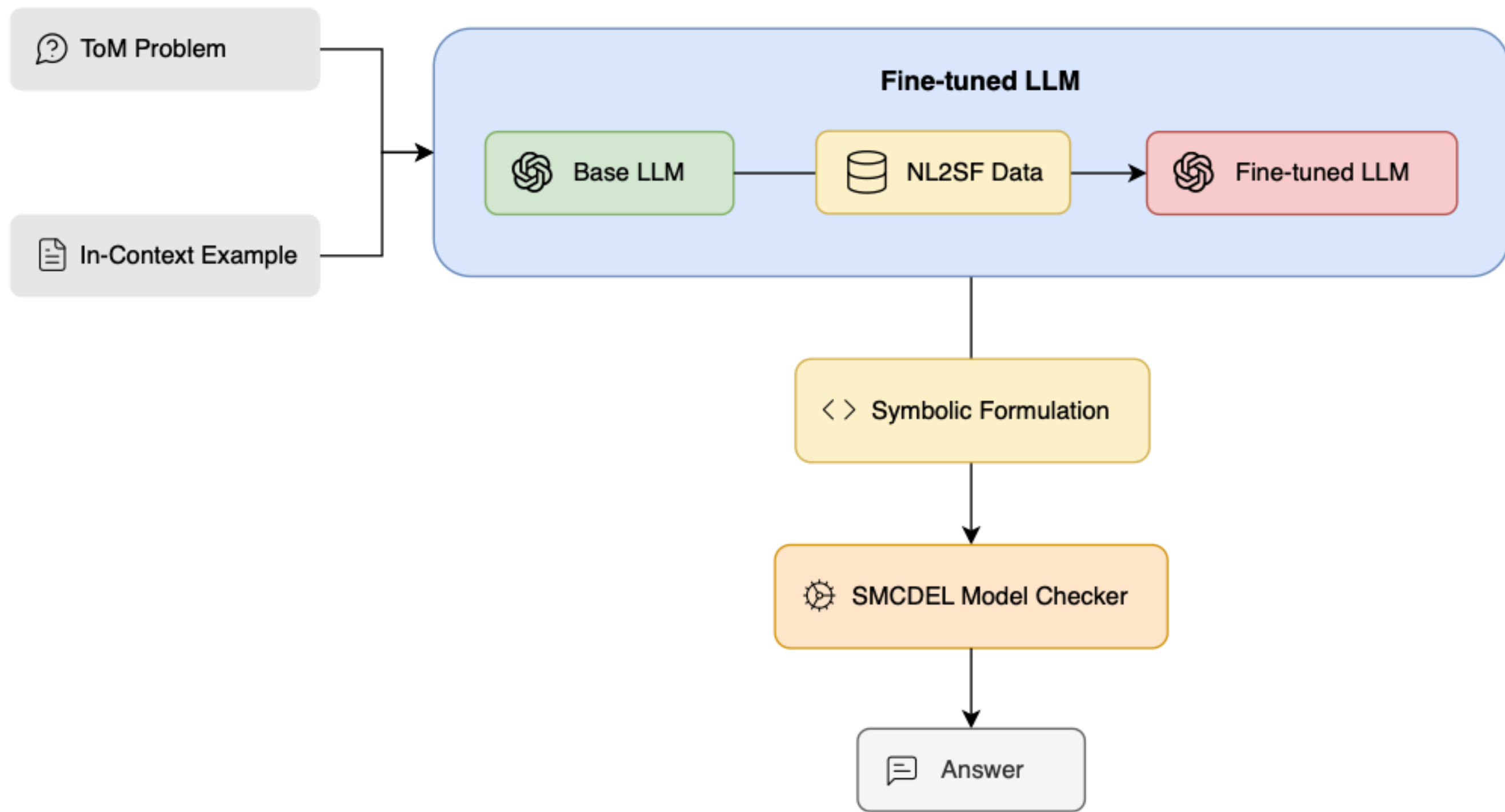
## Symbolic Formulation

VARs 1,2,3,4

LAW Top

OBS Agent<sub>a</sub>:3 Agent<sub>b</sub>:3,4 Agent<sub>c</sub>:1

VALID? [ ! (1|2|3|4) ] [ ! (Agent<sub>a</sub> knows whether (1|2|3|4)) ] (Agent<sub>c</sub> knows whether 2)





# Baseline Slightly better than random guess

Approach	Execution Rate(%)	Accuracy(%)	AUC
DP	99.50	58.00	0.58
SFGP	78.00	49.00	0.60
DP <sub>FT</sub>	<b>100</b>	76.00	0.76
ToM-LM	94.50	<b>91.00</b>	<b>0.94</b>

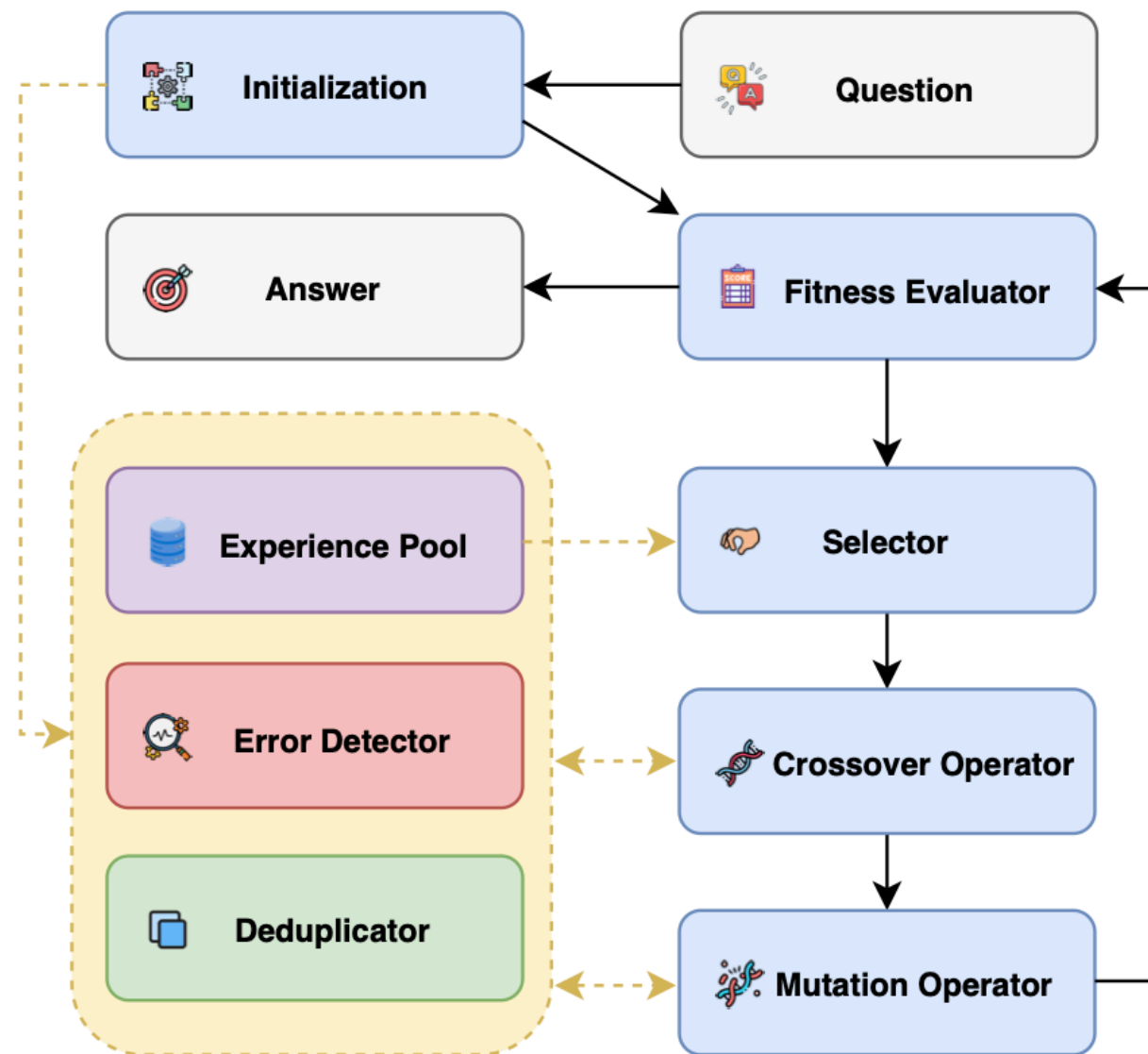
# NP-complete problems

Consider: **Sudoku, Graph Colouring, TSP**

LLMs struggle with multi-objective optimisation,  
strict constraints, and huge search spaces

*Proposal:* Genetic Algorithm meta-framework

Verifier-based vs  
LLM-based modules for  
*Error Detection, Fitness  
Evaluation, Crossover, and  
Mutation*

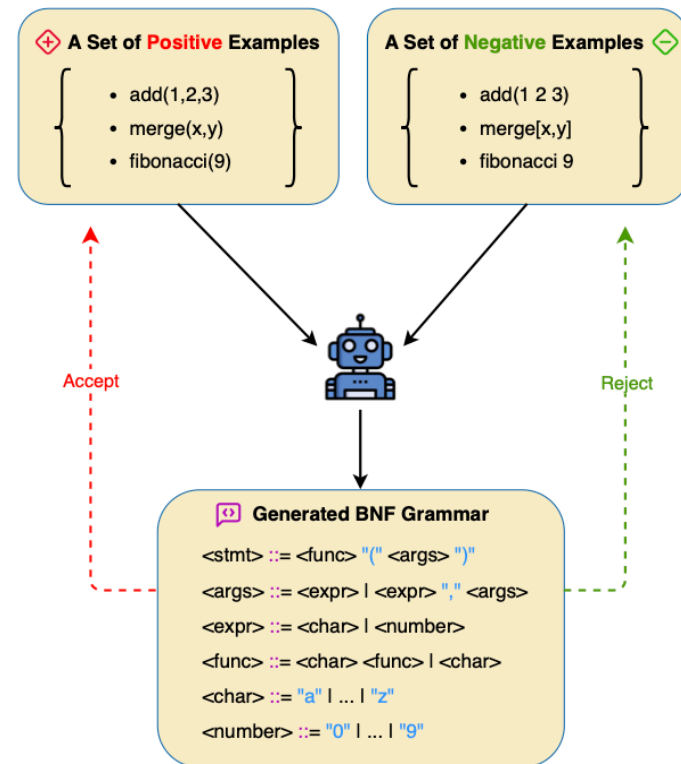


Model	Method	SK <sub>CR</sub>	SK <sub>PS</sub>	GC <sub>CR</sub>	GC <sub>PS</sub>	TSP <sub>CR</sub>	TSP <sub>PS</sub>
GPT-4o-Mini	DP	0	39	0	73	0	79
	BoN	6 $\uparrow 6$	73 $\uparrow 34$	0	86 $\uparrow 13$	4 $\uparrow 4$	94 $\uparrow 15$
	Lyria	<b>8</b> $\uparrow 8$ $\uparrow 2$	<b>73</b> $\uparrow 34$ $-$	0	<b>97</b> $\uparrow 24$ $\uparrow 11$	<b>6</b> $\uparrow 6$ $\uparrow 2$	<b>96</b> $\uparrow 17$ $\uparrow 2$
Qwen2.5:32B-Instruct	DP	0	31	0	74	0	81
	BoN	8 $\uparrow 8$	76 $\uparrow 45$	0	87 $\uparrow 13$	8 $\uparrow 8$	97 $\uparrow 16$
	Lyria	<b>32</b> $\uparrow 32$ $\uparrow 24$	<b>87</b> $\uparrow 56$ $\uparrow 11$	0	<b>96</b> $\uparrow 22$ $\uparrow 9$	<b>30</b> $\uparrow 30$ $\uparrow 22$	<b>99</b> $\uparrow 18$ $\uparrow 2$
Mistral:7B-Instruct	DP	0	0	0	0	0	60
	BoN	0	5 $\uparrow 5$	0	84 $\uparrow 84$	0	80 $\uparrow 20$
	Lyria	0	<b>12</b> $\uparrow 12$ $\uparrow 7$	0	<b>92</b> $\uparrow 92$ $\uparrow 8$	0	<b>89</b> $\uparrow 29$ $\uparrow 9$
Qwen2.5:7B-Instruct	DP	0	26	0	73	0	34
	BoN	0	55 $\uparrow 29$	0	84 $\uparrow 11$	0	88 $\uparrow 54$
	Lyria	0	<b>61</b> $\uparrow 35$ $\uparrow 6$	0	<b>95</b> $\uparrow 22$ $\uparrow 11$	<b>4</b> $\uparrow 4$ $\uparrow 4$	<b>95</b> $\uparrow 61$ $\uparrow 7$

Lyria consistently improves over DP and BoN across models and tasks

Oracle/symbolic Fitness >> LLM Fitness

# Grammar generation



Prompt LLM to produce  $k$  grammars, keep top  $k/2$  by fitness, cross-over and mutate

Combine parser feedback to produce valid BNF initially

Models	Syntax Correctness ( $SX$ )			Semantics Correctness( $SE$ )		
	DP	OPF	HyGenar	DP	OPF	HyGenar
GPT-4o	93	<b>97</b> $\uparrow 4$	96 $\uparrow 3$ $\downarrow 1$	84	85 $\uparrow 1$	<b>93</b> $\uparrow 9$ $\uparrow 8$
GPT-3.5-Turbo	94	95 $\uparrow 1$	<b>99</b> $\uparrow 5$ $\uparrow 4$	37	38 $\uparrow 1$	<b>61</b> $\uparrow 24$ $\uparrow 23$
Llama3:70b-Instruct	57	61 $\uparrow 4$	<b>75</b> $\uparrow 18$ $\uparrow 14$	41	42 $\uparrow 1$	<b>61</b> $\uparrow 20$ $\uparrow 19$
Qwen:72b-Chat	47	49 $\uparrow 2$	<b>76</b> $\uparrow 29$ $\uparrow 27$	20	21 $\uparrow 1$	<b>38</b> $\uparrow 18$ $\uparrow 17$
Mistral:7b-Instruct	1	<b>19</b> $\uparrow 18$	1 $-$ $\downarrow 18$	0	<b>8</b> $\uparrow 8$	1 $\uparrow 1$ $\downarrow 7$
Gemma2:27b-Instruct	91	92 $\uparrow 1$	<b>98</b> $\uparrow 7$ $\uparrow 6$	56	57 $\uparrow 1$	<b>79</b> $\uparrow 23$ $\uparrow 22$
Starcoder2:15b-Instruct	76	60 $\downarrow 16$	<b>98</b> $\uparrow 22$ $\uparrow 38$	30	20 $\downarrow 10$	<b>44</b> $\uparrow 14$ $\uparrow 24$
Codestral:22b	53	71 $\uparrow 18$	<b>80</b> $\uparrow 27$ $\uparrow 9$	44	52 $\uparrow 8$	<b>67</b> $\uparrow 23$ $\uparrow 15$

# Takeaways

Thinking "**fast and slow**"

Integration: using formal systems to **vet** LLM solutions and as "*Sources of Truth*"

*Other successes:* AlphaGeometry and AlphaProof demonstrate viability of hybrid approaches