



Systems Opportunities for LLM Fine-Tuning using Reinforcement Learning

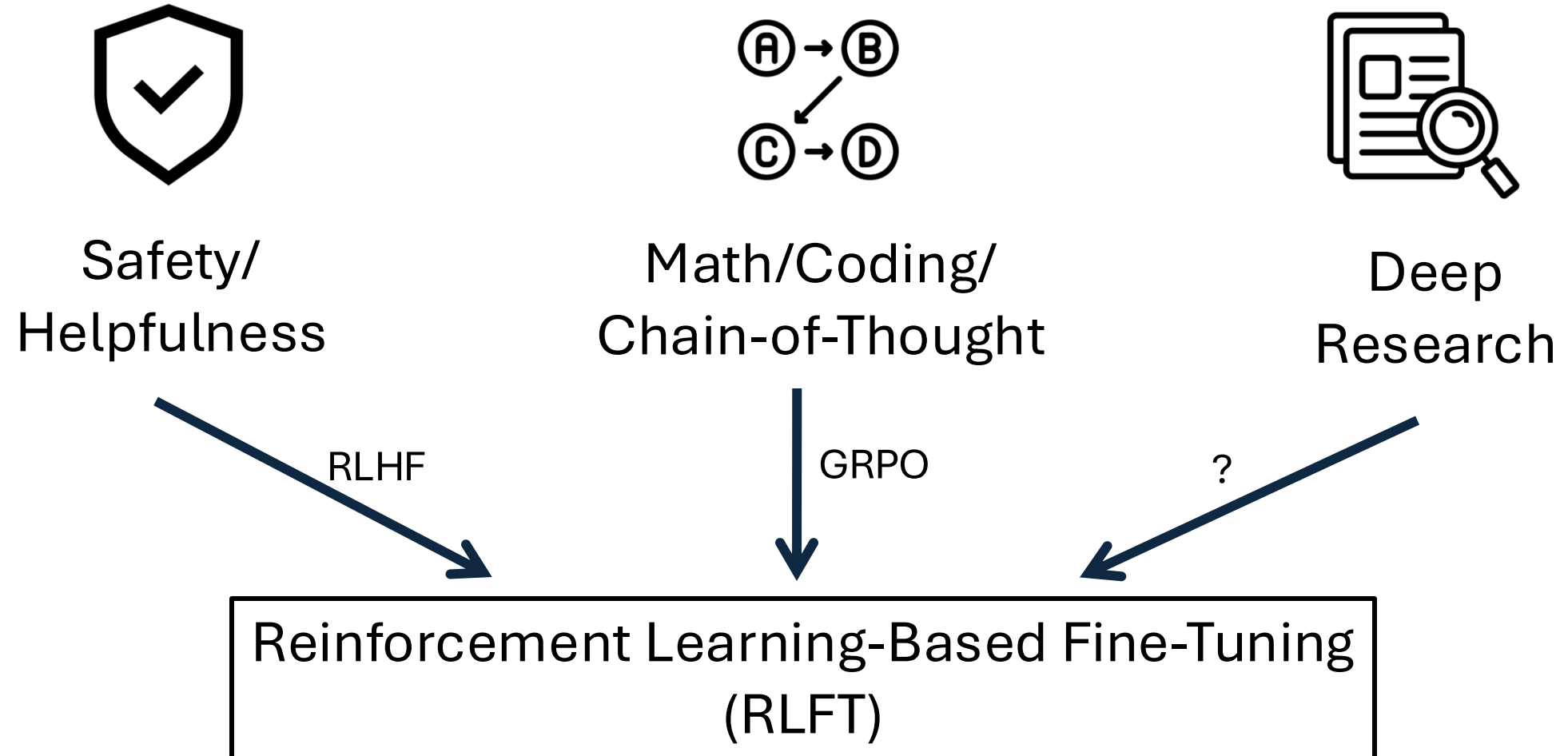
Pedro F. Silvestre, Peter Pietzuch

pms20@ic.ac.uk

Imperial College London

Large Scale Data & Systems (LSDS) Group

Recent Explosion in LLM Capabilities

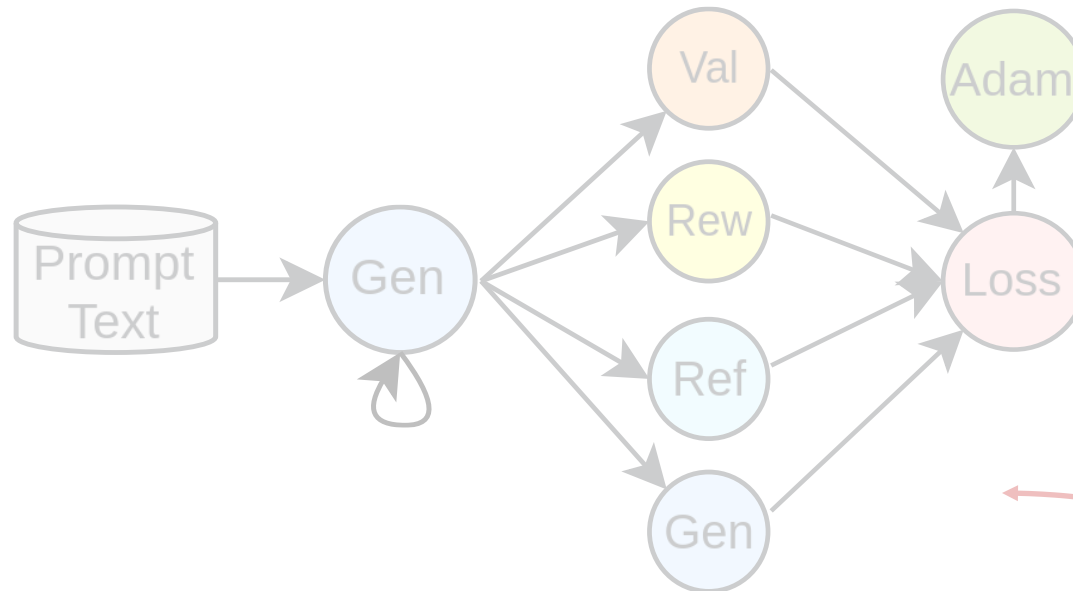


RLFT Generates Variable-Size Data

Context



← **Variable-Size**



← **Large-Scale**

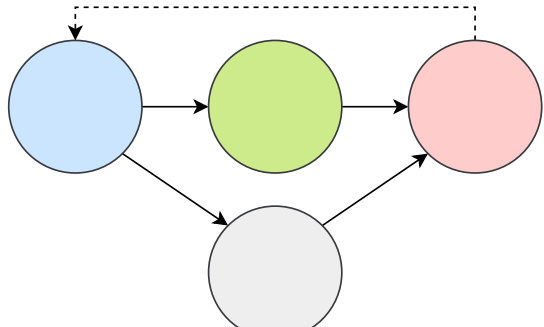
1. Sequential
Generation

2. Parallel
Inference

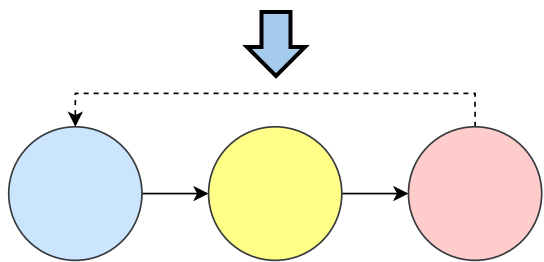
3. Backward
+ Training

Tension Between **Large-Scale** and **Variable Size**

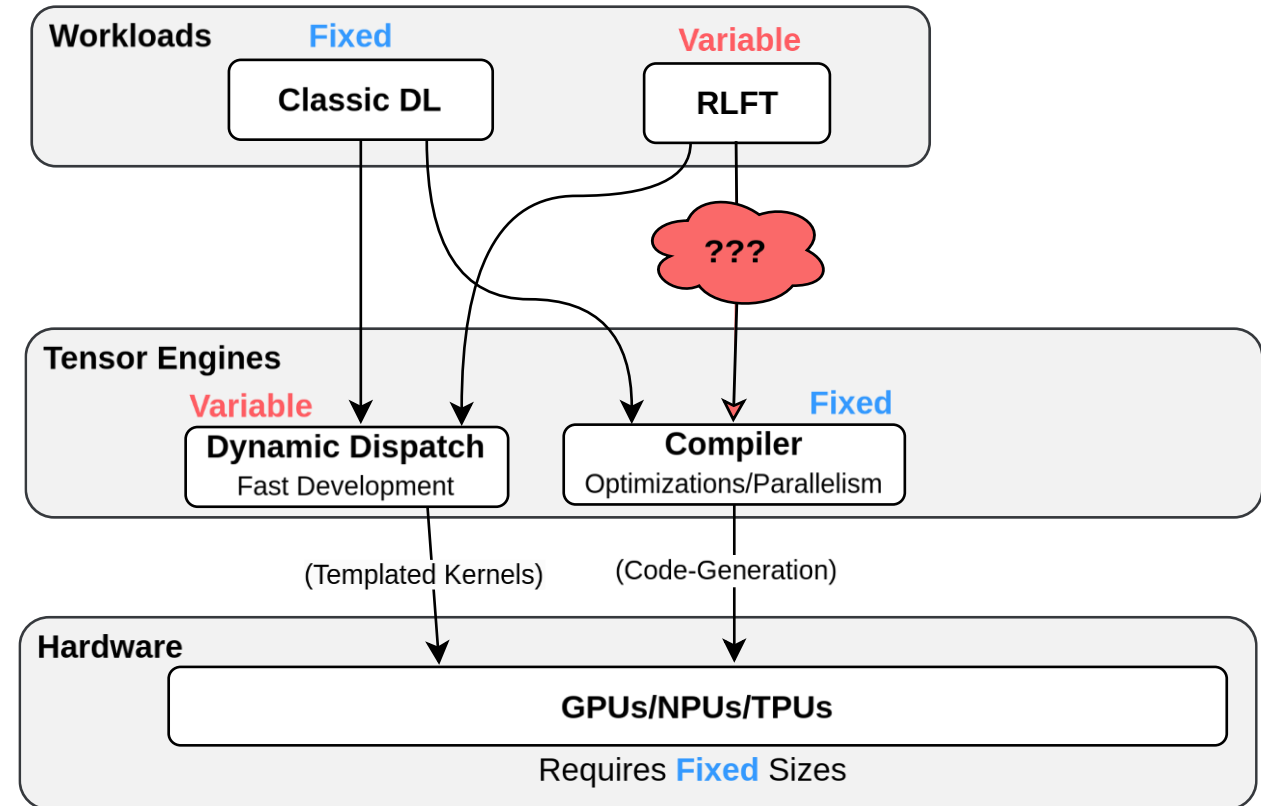
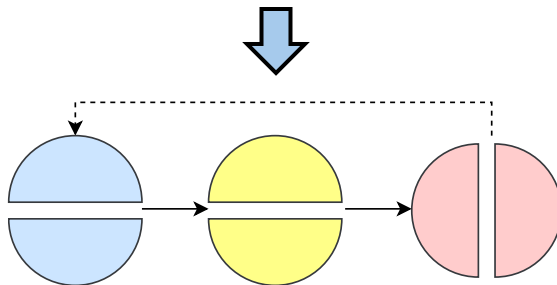
Computation Graph



Optimizations



**Parallelization/
Distribution**



Workarounds for Mapping RLFT To Compilers

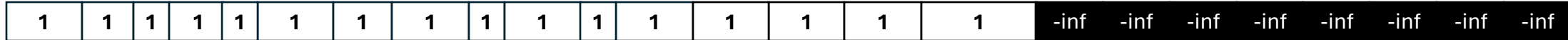
Undershoot: Error

Overshoot: Overhead

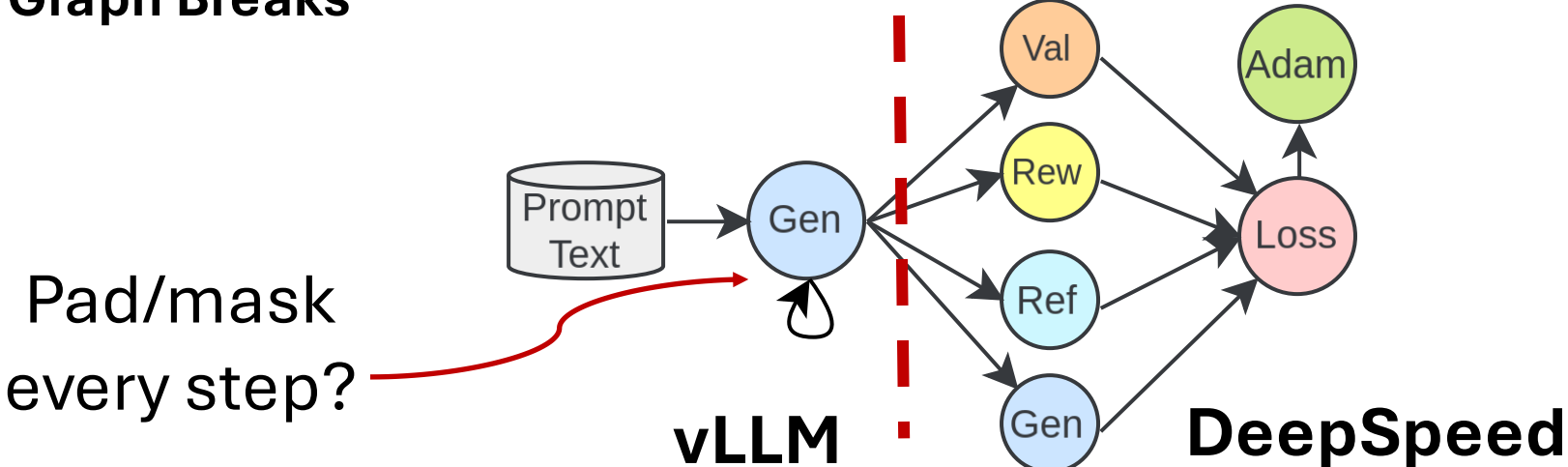
1. Pad to fixed MAX_LEN



2. Mask valid steps



3. Graph Breaks



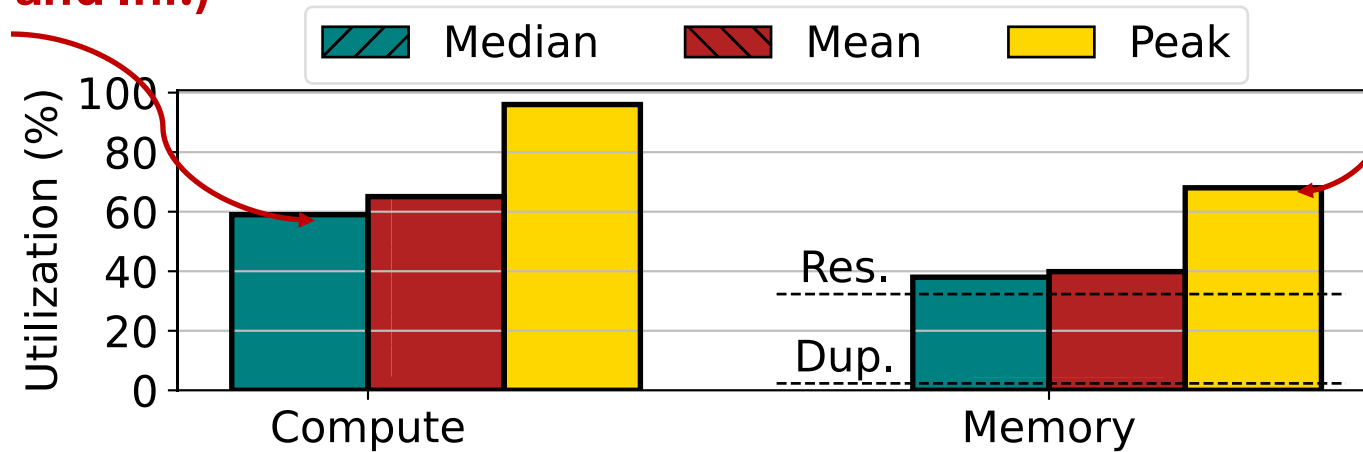
Manual Work!

Not Holistic!

TRL GRPO Benchmark on A100

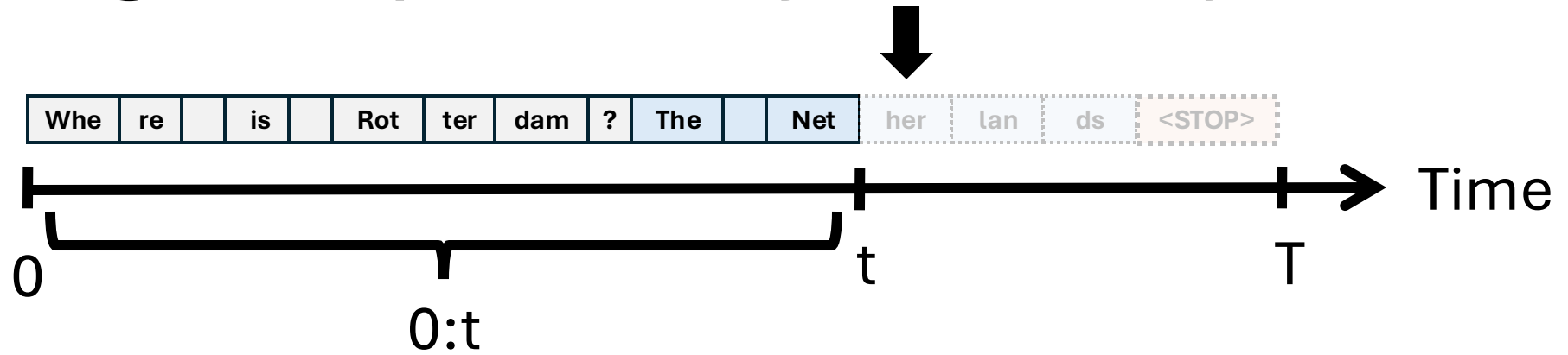
GPU underutilization
(No Overlap of Gen. and Inf.)

BS=2 exceeds memory
(Large Peak Memory)

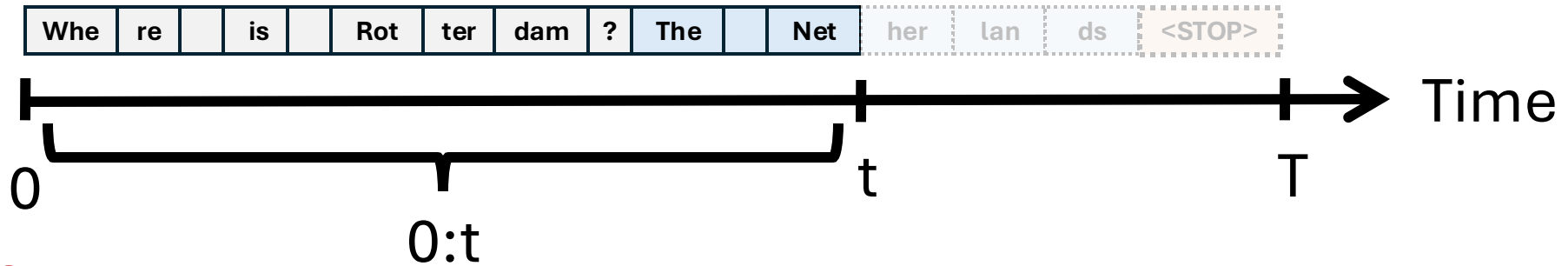


Need Better Graph Representations for RLFT!

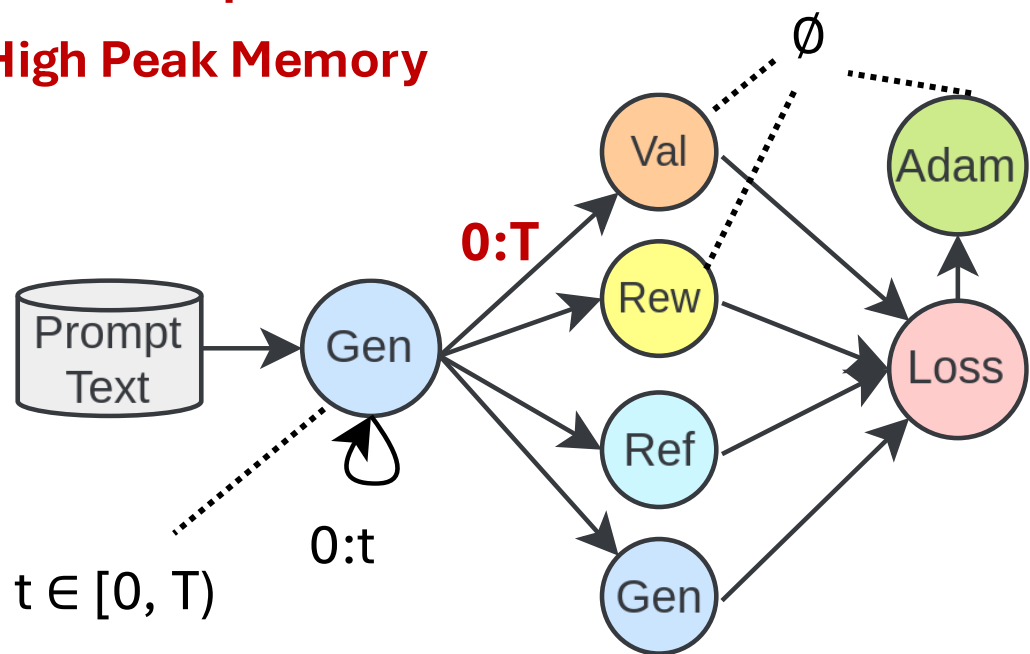
Key Insight: Represent Symbolically



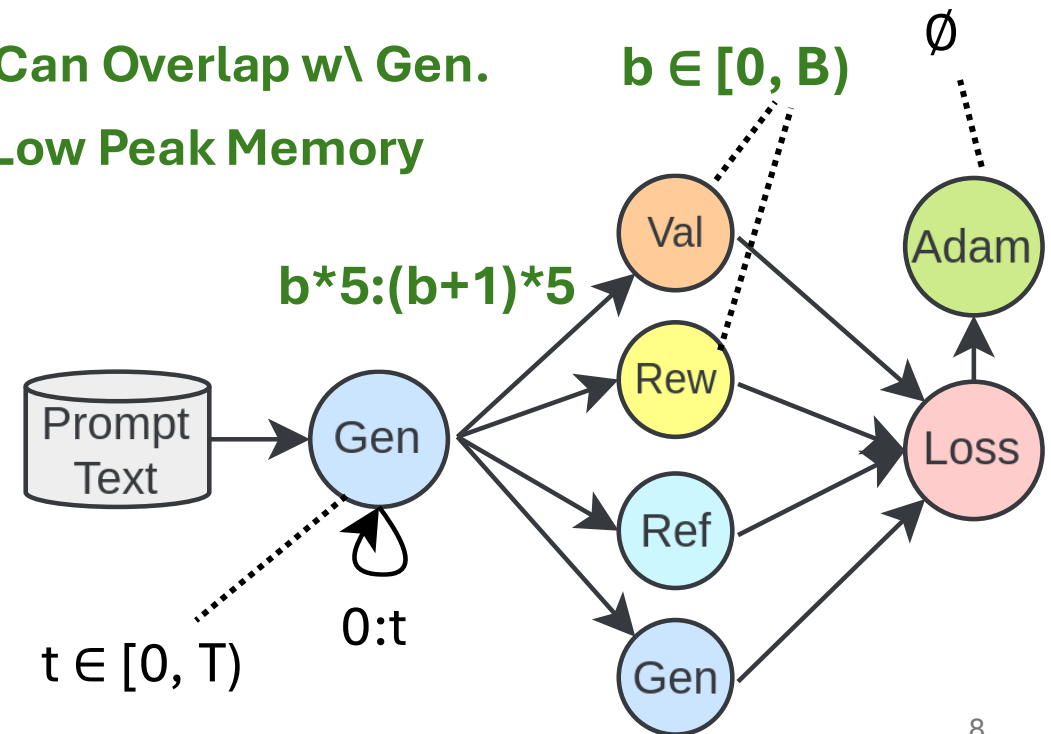
Key Insight: Represent Symbolically



No Overlap w/ Gen.
High Peak Memory



Can Overlap w/ Gen.
Low Peak Memory



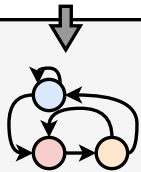
$B = T/5$

TimeRL: A Symbolic Deep Learning Compiler

Expose Symbols in Programming Model

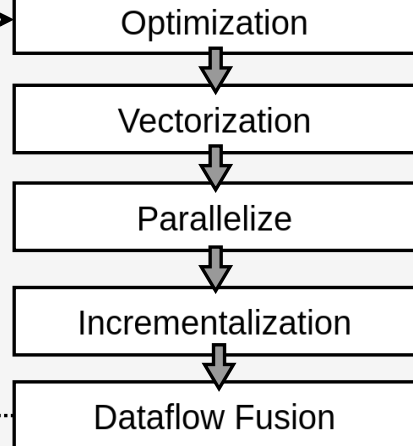
(No pads/masks!)

```
ctx = trl.like()
ctx[0:P] = prompts.get()
ctx[t] = gen(ctx[0:t])
r = rew(ctx[0:T])
ref_log = ref(ctx[0:T])
loss = loss_fn(ctx, ref_log, r)
loss.backward()
```



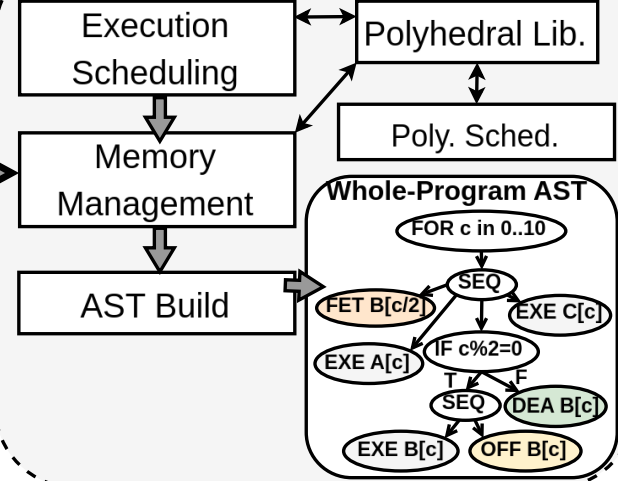
Holistic Graph Transformations

(Min. Variable-shapes)



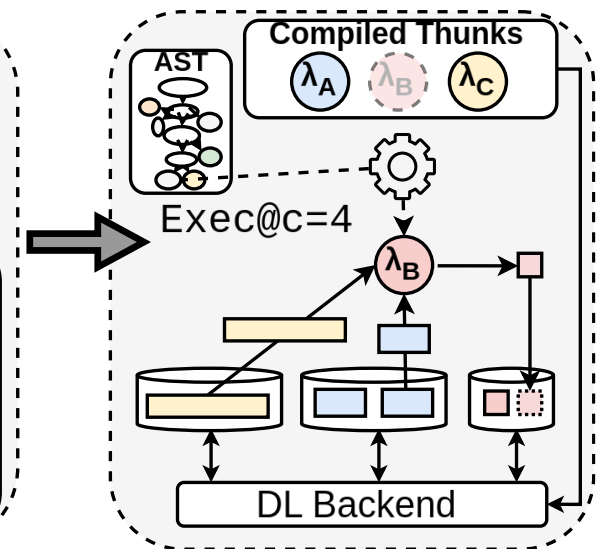
Whole-Program Scheduling

(Emit a program AST)



Evaluate Symbols at Runtime

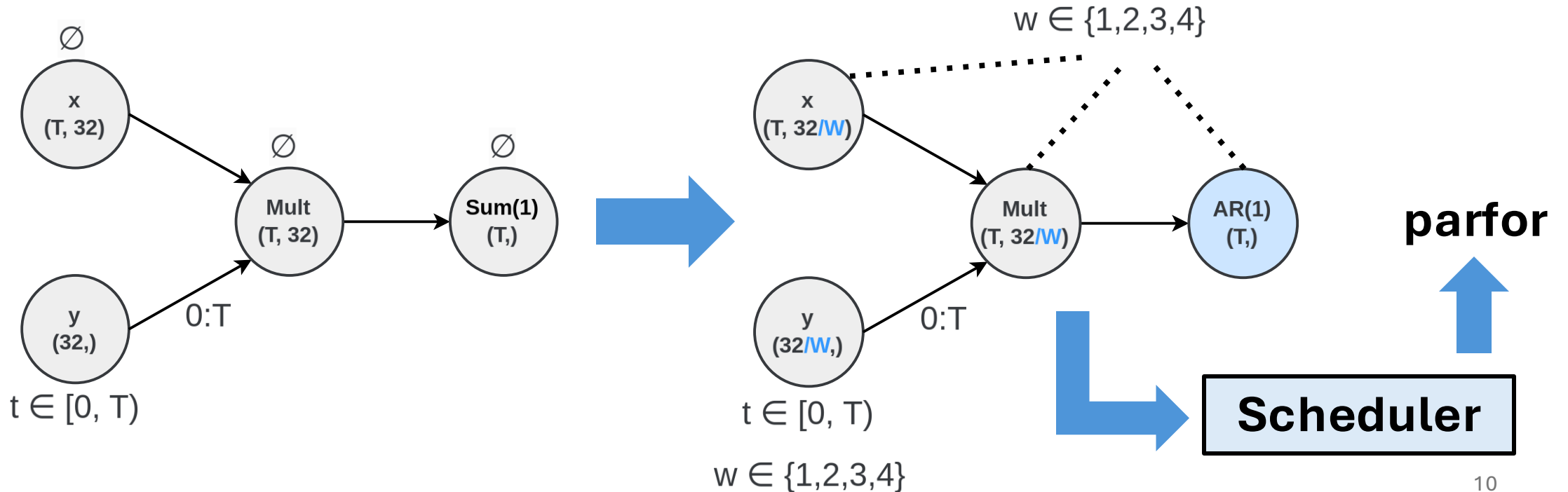
(Specialized Storage)



How Symbols Simplify Parallelism

- Intra-Operator (TP, DP, SP):
 - **CollectiveOps**
 - **Worker** dimension

- Inter-Operator (PP, EP):
 - **P2POps**
 - **Microbatch** dimension



Come Find Out More!

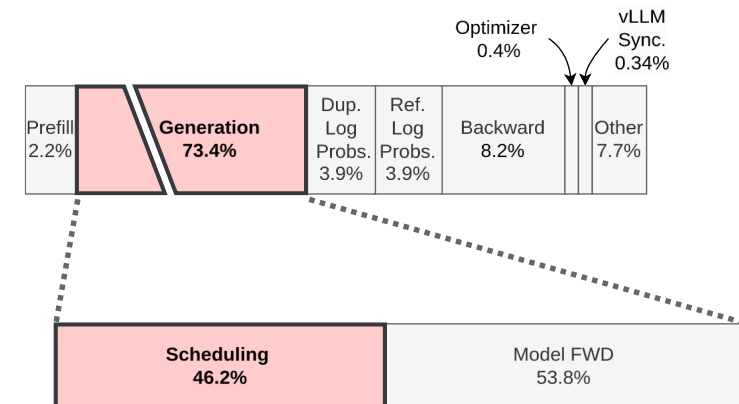
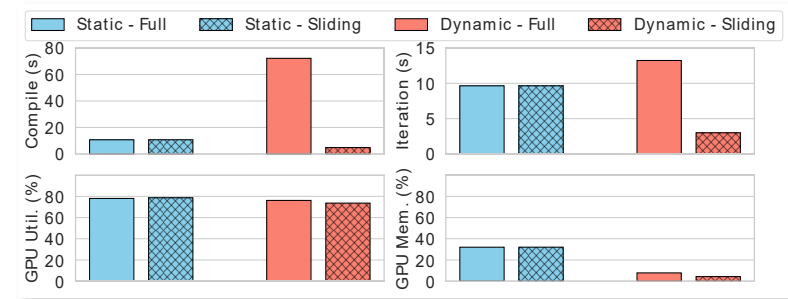
- Challenges & Opportunities
 - Survey of existing RLFT Systems

- Evaluation of JAX for Generation
 - Overhead due to padding/masking.
 - Especially for window attention.

- Discussion on vLLM for Generation
 - Scheduling overheads!
 - Propose a method to address this

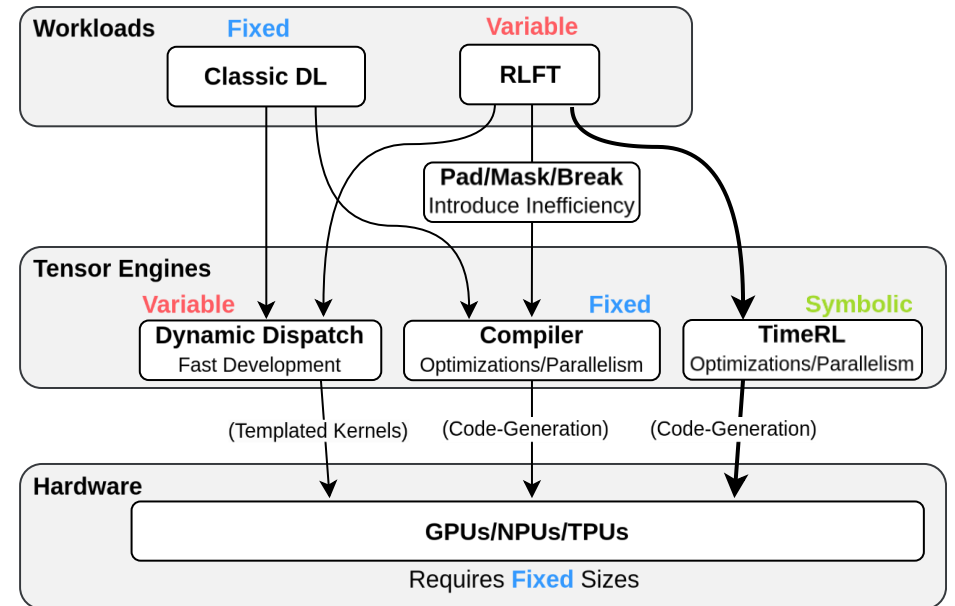
System	DL Eng.	Graph Based	Multi Algo.	Unif. Eng.	Par. Space
NeMo-Aligner[55]	PT	✗	✓	✓	3D [§]
OpenRLHF[26]	PT	✗	~	✗*	3D [§]
DS-Chat[71]	PT	✗	~	✓	3D [§]
TRL[65]	PT	✗	✓	✗*	3D [§]
FlexRLHF[69]	PT	✗	✗	✓	3D [§]
Puzzle[36]	PT	✗	✗	✗‡	3D [§]
ReaLHF[44]	PT	✓ [◇]	✓	✓	3D
RLHFuse[79]	PT	✓ [◇]	✓	✗ [†]	3D
HybridFlow[56]	PT	✓ [◇]	✓	✗*	3D

* vLLM for generation. † In-house vLLM-like engine for generation.
 ‡ DeepSpeed-Inference for generation. § Fixed strategy. ◇ Coarse-grained.

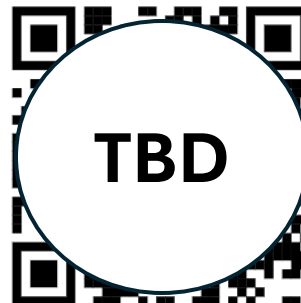


Conclusion

- RLFT generates **variable**-size contexts
- Existing compilers need **fixed** sizes
- Workarounds are **manual and not holistic**
- **TimeRL**: **symbols** to represent variables



TimeRL



RLFT Opportunities