# POET: Training Neural Networks on Tiny Devices with Integrated Rematerialization and Paging

## Shishir G. Patil

With Paras Jain, Prabal Dutta, Ion Stoica, Joseph Gonzalez
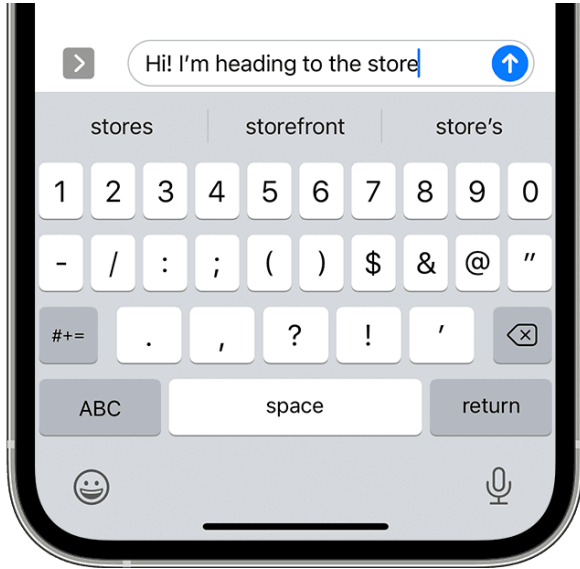
https://github.com/ShishirPatil/poet

ICML
International Conference
On Machine Learning

riselab
UC Berkeley

BAIR
BERKELEY ARTIFICIAL INTELLIGENCE RESEARCH

# POET: Training Neural Networks on Tiny Devices with Integrated Rematerialization and Paging

Shishir G. Patil

BERT on edge devices!

With Paras Jain, Prabal Dutta, Ion Stoica, Joseph Gonzalez

https://github.com/ShishirPatil/poet

ICML International Conference On Machine Learning

riselab UC Berkeley

BAIR BERKELEY ARTIFICIAL INTELLIGENCE RESEARCH

# Model Personalization Adapts Models by Training on User Data to Improve Accuracy



**Autocompletion**

**Voice Recognition**

**Fitness Tracker**

# Model Fine-tuning – Train on Edge

**Fine-tune on-device**

Train

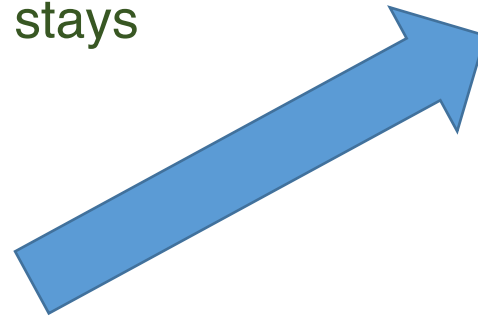**Key Challenge: Limited memory for DNN training!**

**Pros:**
+ guarantees user's privacy as all data stays on their device
+ enables offline device operation

**Cons:**
- cannot train modern DNNs on edge devices
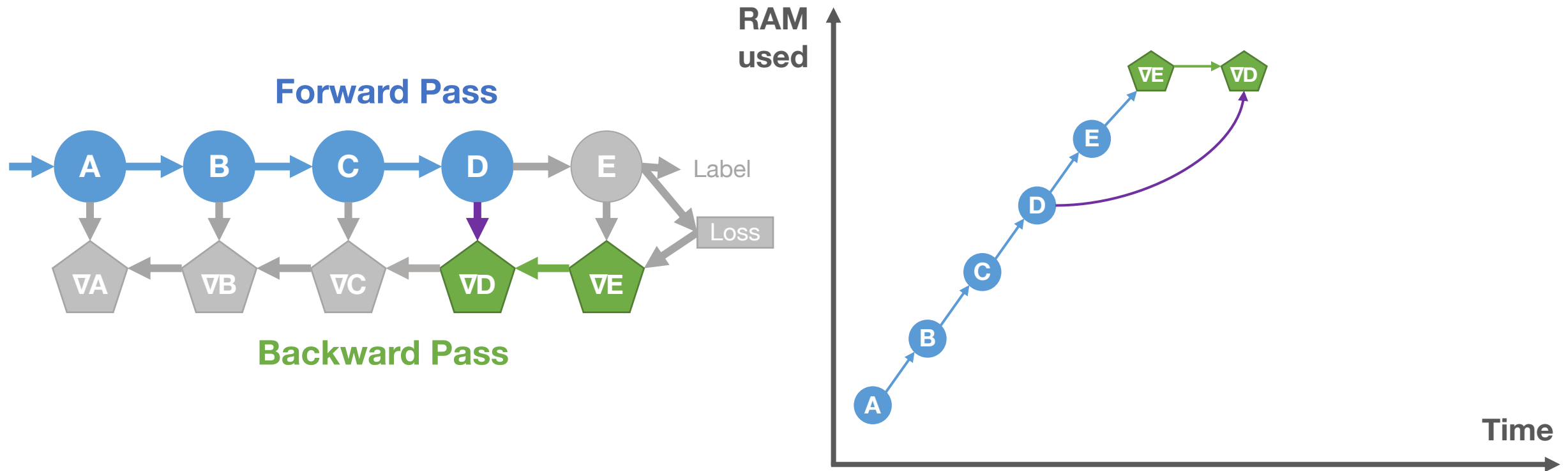
How to reduce the memory and energy requirements of ML training for modern DNN architectures within the constraints of edge devices?

# Training is Memory Intensive since Activation from Forward Pass Need to be Stored for Backpropagation

**Forward Pass**

A → B → C → D → E → Label

Loss

∇A ← ∇B ← ∇C ← ∇D ← ∇E

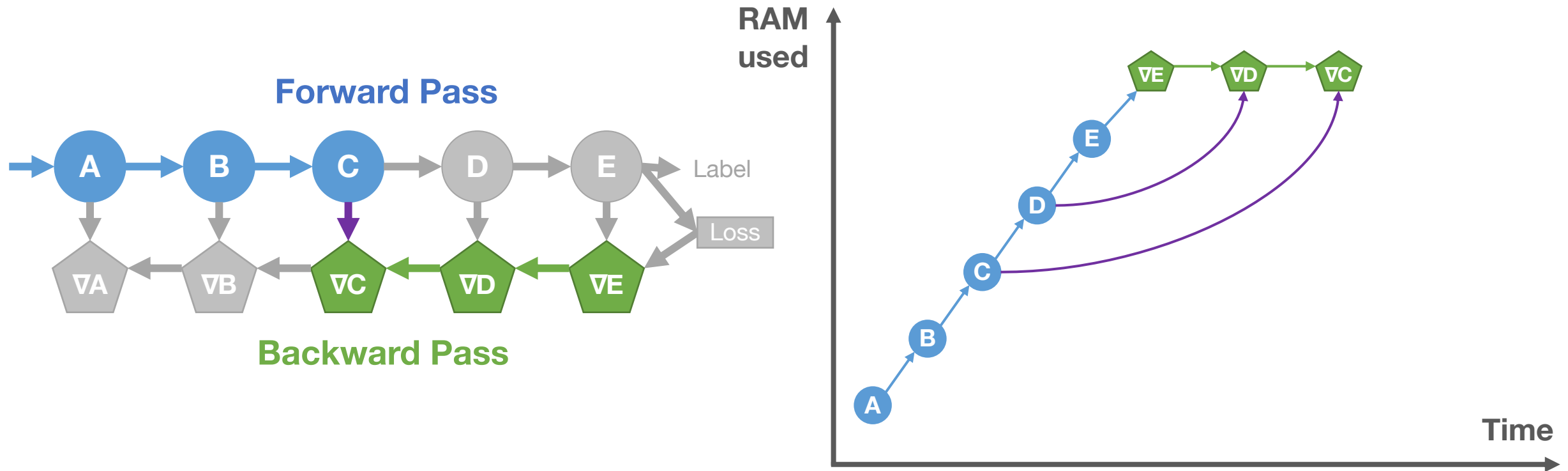**Backward Pass**

RAM used

∇E

E

D

C

B

A

Time
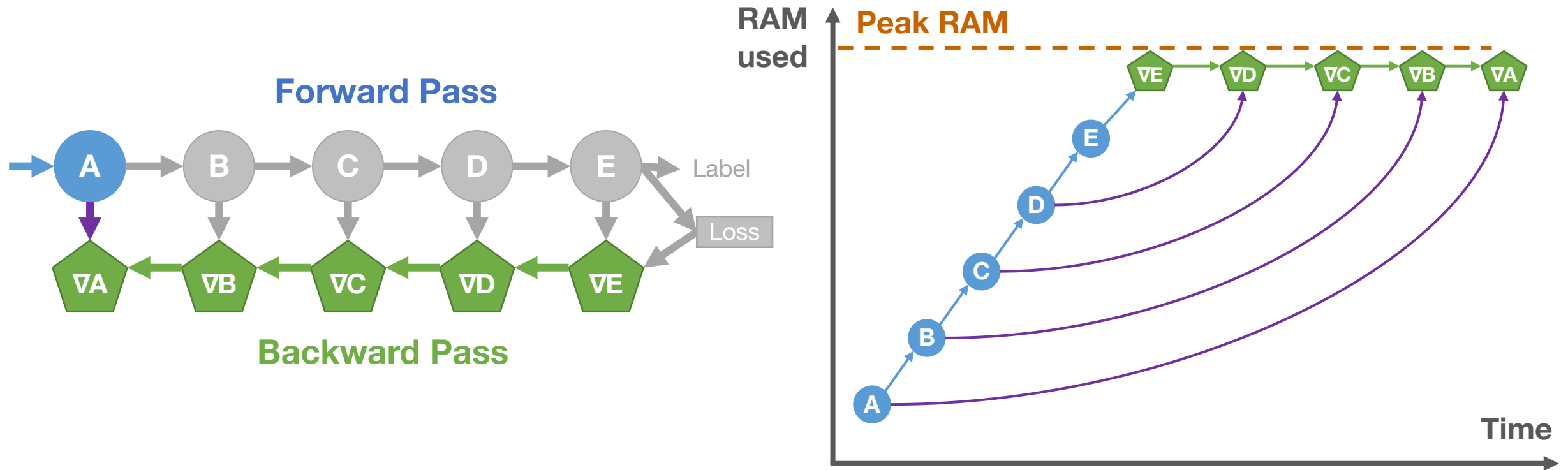
# Training is Memory Intensive since Activation from Forward Pass Need to be Stored for Backpropagation

# Training is Memory Intensive since Activation from Forward Pass Need to be Stored for Backpropagation

# Training is Memory Intensive since Activation from Forward Pass Need to be Stored for Backpropagation



**Forward Pass**

**Backward Pass**

Label

Loss

RAM used

**Peak RAM**

Time

# Rematerialization and Paging: Two Techniques to Reduce Memory Consumption



**Forward Pass**

**Backward Pass**

Label

Loss

**RAM used**

Peak RAM (no rematerialization nor paging)

Time

# Rematerialization:
Free early & recompute

# Rematerialization and Paging: Two Techniques to Reduce Memory Consumption



**Forward Pass**

**Backward Pass**

Label

Loss

**RAM used**

Peak RAM (no rematerialization nor paging)

**Time**

# Rematerialization:
Free early & recompute

# Rematerialization and Paging: Two Techniques to Reduce Memory Consumption



**Forward Pass**

A → B → C → D → E → Label

Loss

∇A ← ∇B ← ∇C ← ∇D ← ∇E

**Backward Pass**

RAM used

Peak RAM (no rematerialization nor paging)

Available RAM

Peak RAM

Time

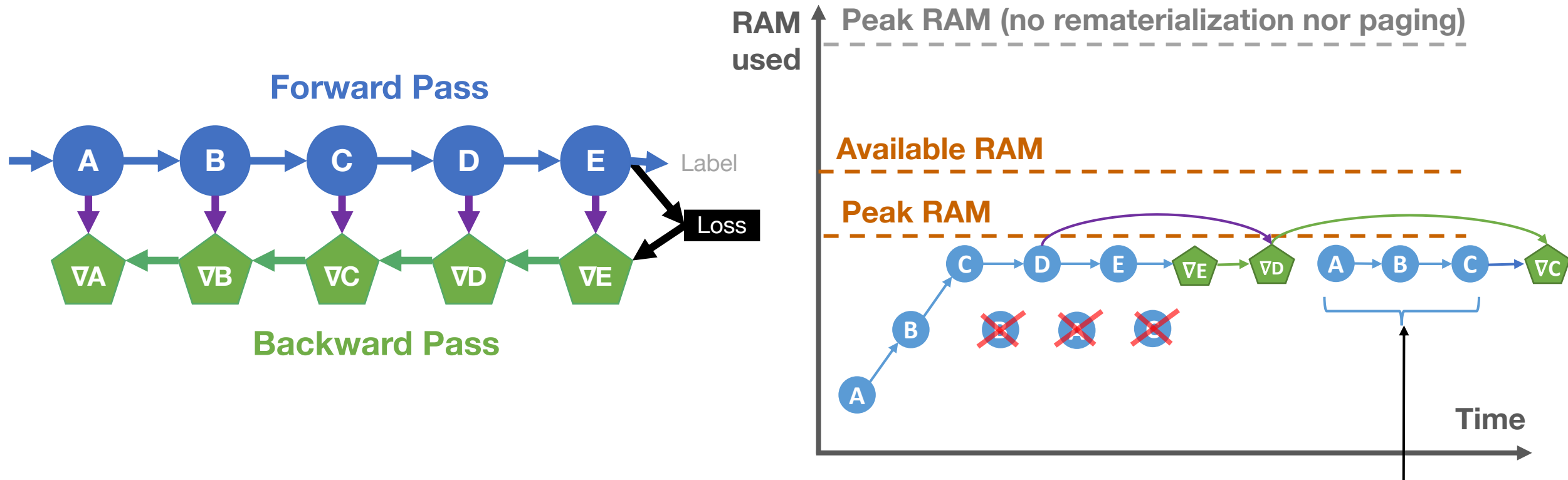Additional Energy and runtime due to recomputation!

# Rematerialization:
Free early & recompute

# Rematerialization and Paging: Two Techniques to Reduce Memory Consumption

**Forward Pass**

**Backward Pass**

Label

Loss

RAM used

Peak RAM (no rematerialization nor paging)

**Available RAM**

**Peak RAM**

**Time**

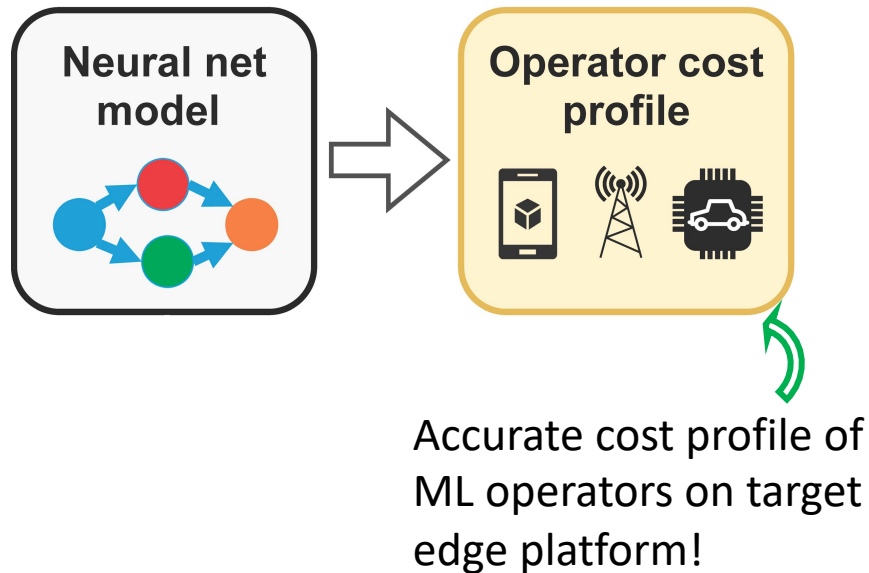Additional Energy due to Paging

**Paging:**
Page-out to secondary storage and page-in Just-in-Time!

# POET: Private Optimal Energy Training



Neural net model

{SOTA ML model, memory and runtime constraints}

# POET: Private Optimal Energy Training



Accurate cost profile of
ML operators on target
edge platform!

**POET: Training Neural Networks on Tiny Devices**     https://poet.cs.berkeley.edu

# POET: Private Optimal Energy Training



**Neural net model** → **Operator cost profile** → **POET solver**

POET solver
min *total energy usage*
s.t. *memory constraint*
s.t. *runtime constraint*

Incorporate *memory* and *runtime* constraints into a Mixed Integer Linear Program (MILP) formulation

**POET:  Training Neural Networks on Tiny Devices**                    https://poet.cs.berkeley.edu

# POET: Private Optimal Energy Training



**Neural net model** → **Operator cost profile** → **POET solver**

min *total energy usage*
s.t. *memory constraint*
s.t. *runtime constraint*

→ **Execute on edge device**
(1) Rematerialize ✓→✗
(2) Page to flash ⬛→💾

POET finds a provably optimal solution through integrated rematerialization and paging.

# POET: Private Optimal Energy Training

Layer of Neural Network



Logical Time

Forward

Backward

Pixelated box indicates activation tensor for Layer '*l*' is resident in RAM at timestep '*t*'

POET's integrated rematerialization and paging search space finds advanced solutions that are not possible through simple heuristics.

**POET: Training Neural Networks on Tiny Devices**

https://poet.cs.berkeley.edu

# Result: POET lowers energy consumption and allows training large models previously not possible!



ResNet18 on Cortex A72

Lower is better

| PyTorch baseline | DTR (Kirasame et al. 2021) | Checkmate (Jain et al. 2020) |
| POET (ours) | revolve (Griewank and Walther 2000) | Chen et al. 2016 |

**POET: Training Neural Networks on Tiny Devices**

https://poet.cs.berkeley.edu

# Result: POET lowers energy consumption and allows training large models previously not possible!



ResNet18 on Cortex A72

Relative Energy Usage FWD+BWD

Activation RAM savings

Naïve Strategy (Tensorflow / Torch)

- ★- PyTorch baseline
- --+-- DTR (Kirasame et al. 2021)
- --■-- Checkmate (Jain et al. 2020)
- --•-- POET   (ours)
- --•-- revolve (Griewank and Walther 2000)
- --♦-- Chen et al. 2016

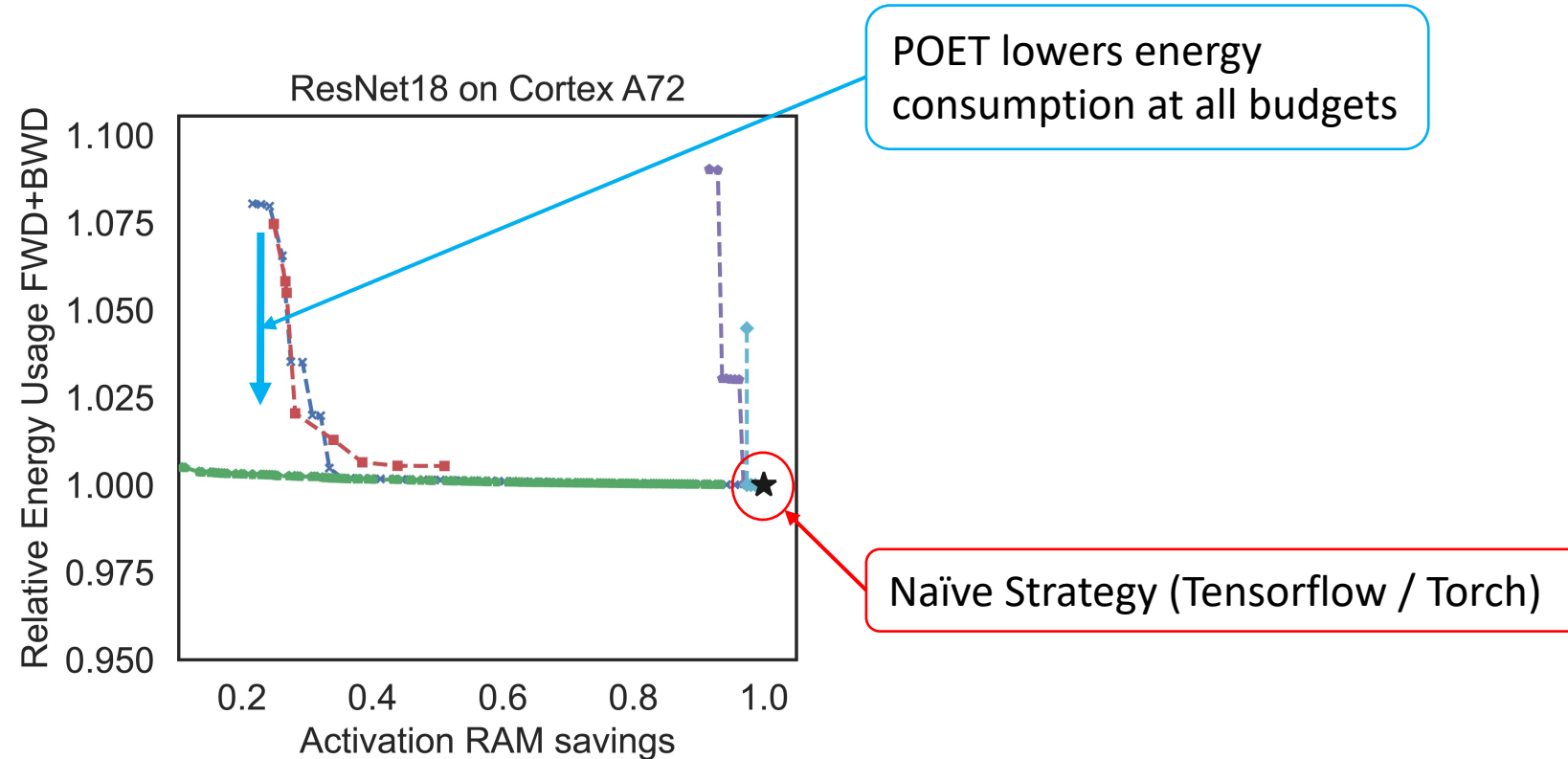# Result: POET lowers energy consumption and allows training large models previously not possible!



POET lowers energy consumption at all budgets

Naïve Strategy (Tensorflow / Torch)

**ResNet18 on Cortex A72**

Legend:
- PyTorch baseline
- POET (ours)
- DTR (Kirasame et al. 2021)
- revolve (Griewank and Walther 2000)
- Checkmate (Jain et al. 2020)
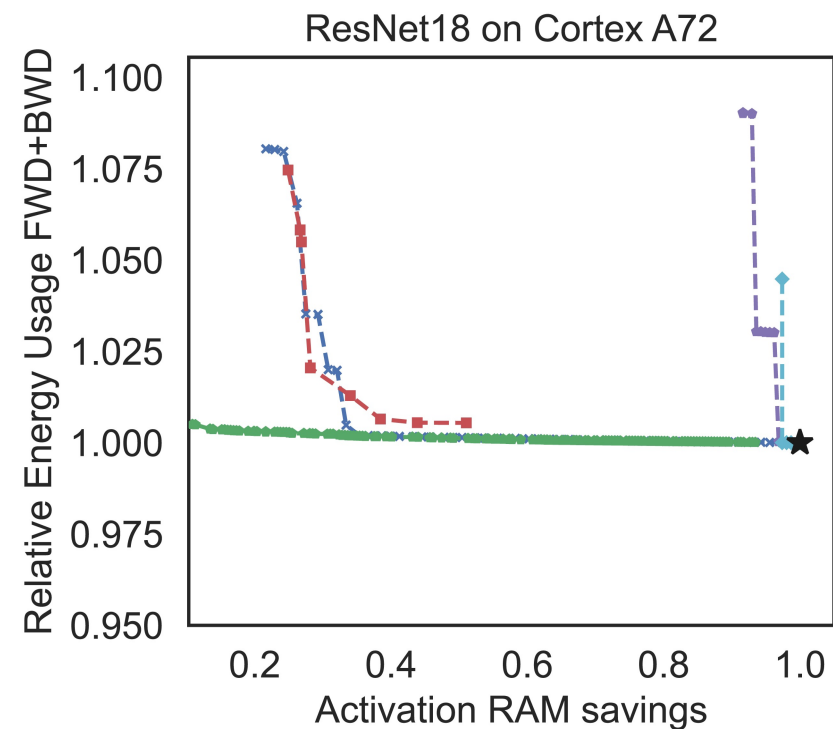- Chen et al. 2016

**POET:  Training Neural Networks on Tiny Devices**

https://poet.cs.berkeley.edu

# Result: POET lowers energy consumption and allows training large models previously not possible!



POET lowers energy consumption at all budgets

POET's integrated Rematerialization and Paging enables training with much smaller memory budgets which was previously not possible!

Naïve Strategy (Tensorflow / Torch)

ResNet18 on Cortex A72

PyTorch baseline
POET (ours)
DTR (Kirasame et al. 2021)
revolve (Griewank and Walther 2000)
Checkmate (Jain et al. 2020)
Chen et al. 2016

**POET: Training Neural Networks on Tiny Devices**

https://poet.cs.berkeley.edu

# Result: POET lowers energy consumption and allows training large models previously not possible!



POET's integrated Rematerialization and Paging enables training with much smaller memory budgets which was previously not possible!

Legend:
- PyTorch baseline
- POET (ours)
- DTR (Kirasame et al. 2021)
- revolve (Griewank and Walther 2000)
- Checkmate (Jain et al. 2020)
- Chen et al. 2016

**POET: Training Neural Networks on Tiny Devices**

https://poet.cs.berkeley.edu

## Conclusion

- POET enables training SOTA DNN models locally on memory-constrained edge devices.

- POET's fine grained profiling results in accurate cost profiles.

- POET's MILP formulation finds the optimal training schedule through integrated **rematerialization** and **paging.**

🌐 https://poet.cs.berkeley.edu

✉ shishirpatil@berkeley.edu

https://github.com/ShishirPatil/poet