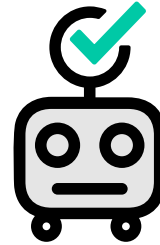# Checkmate

**Checkmate:** Breaking the Memory Wall with Optimal Tensor Rematerialization

Paras Jain

*Joint work with:* Ajay Jain, Ani Nrusimha, Amir Gholami, Pieter Abbeel, Kurt Keutzer, Ion Stoica, Joseph Gonzalez

**BigGAN (2018)**

Image generation

Brock et al. 2019

**VideoBERT (2019)**

Video generation

VideoBERT

Sun et al. 2019

**GPT-2 (2019)**

Text generation

SYSTEM PROMPT (HUMAN-WRITTEN)

*In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.*
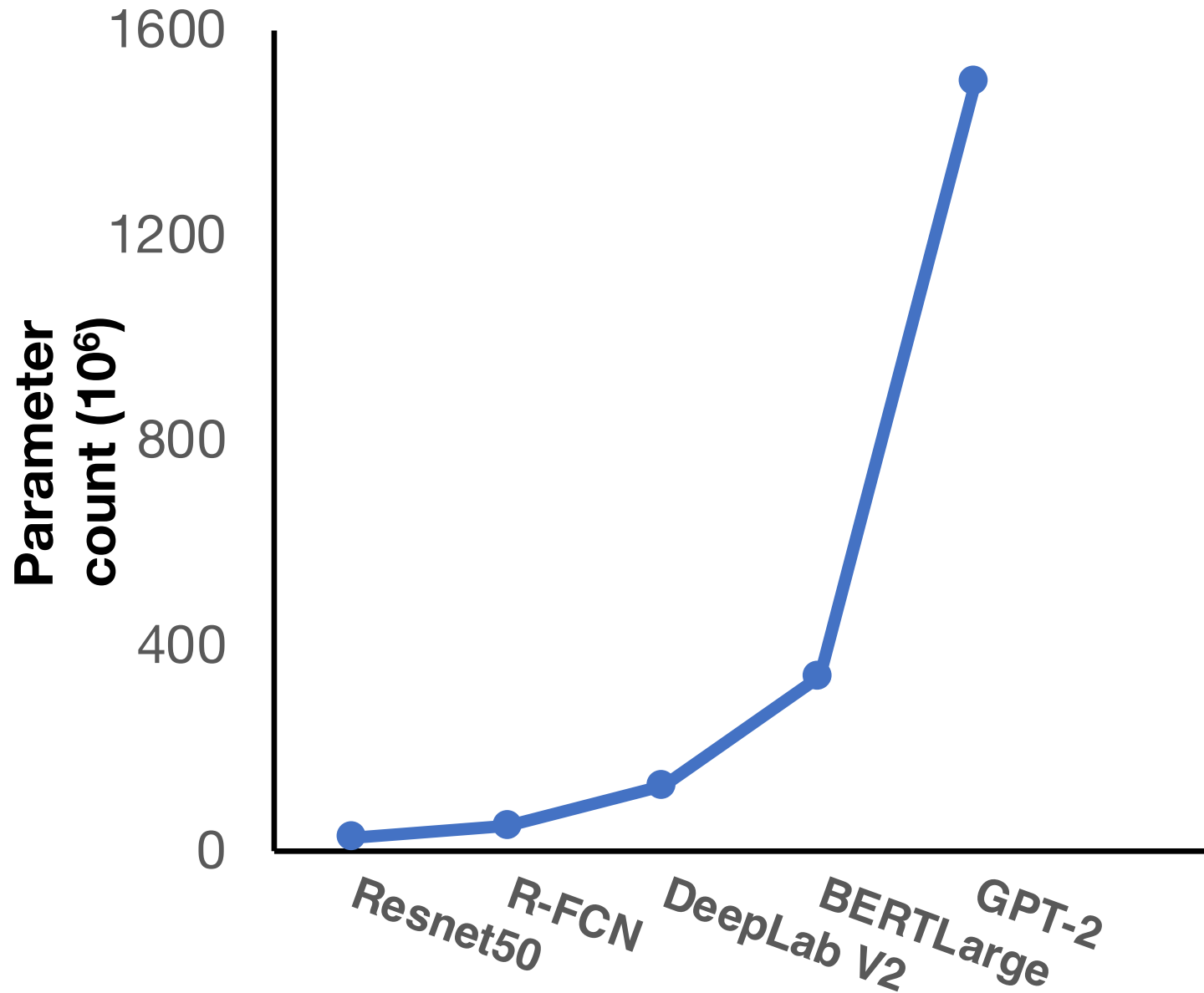
MODEL COMPLETION (MACHINE-WRITTEN, 10 TRIES)

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.
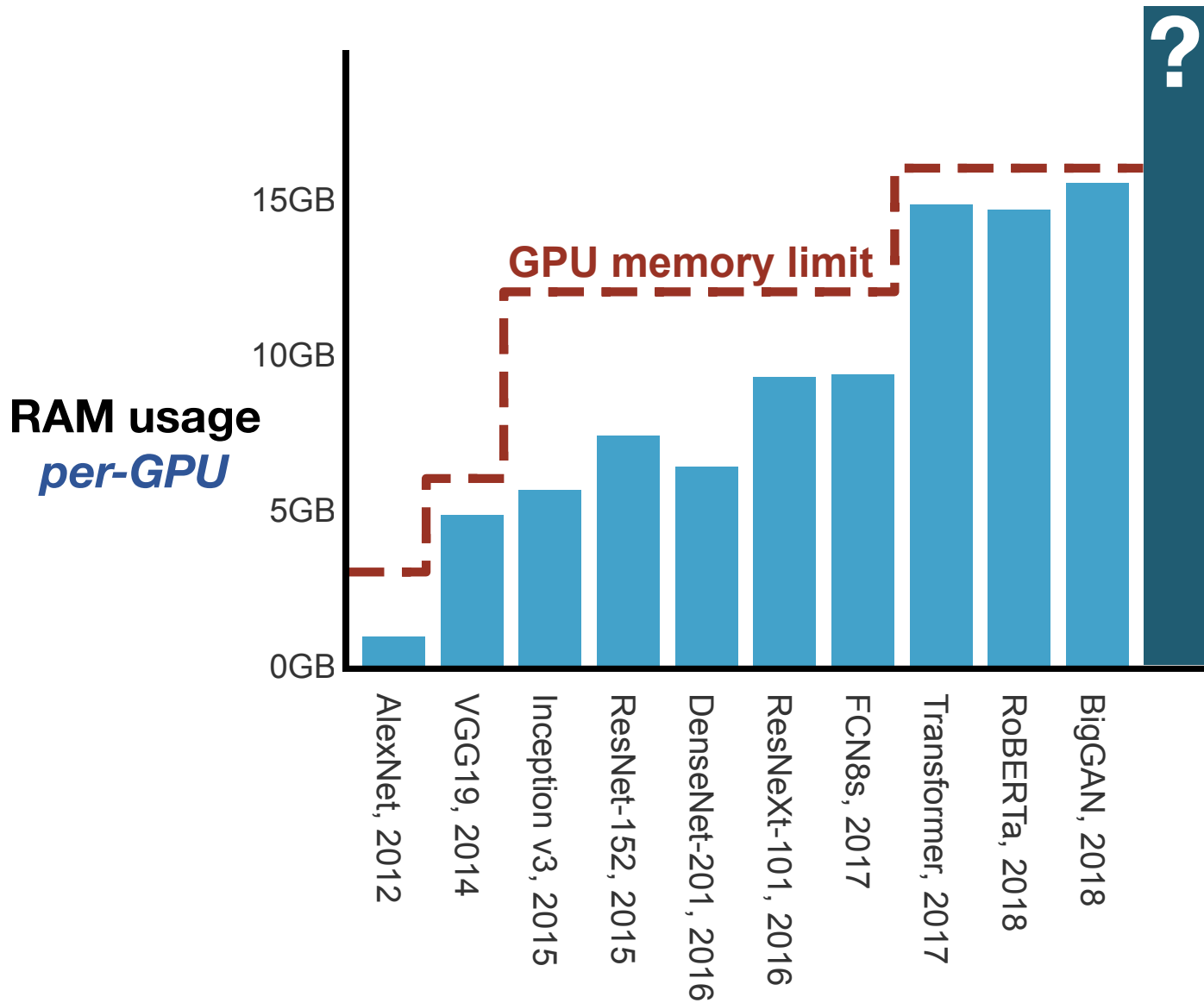
Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several

Radford et al. 2019

**Emerging trend:**

Rapid **growth in model size**

Figure adapted from NVIDIA

checkmateai.github.io

**RAM usage** *per-GPU*

GPU memory limit

State-of-the-art models have hit a **memory capacity wall**.

**Cited memory as limiting factor**

Chen et al. 2016          Liu et al. 2019
Gomez et al. 2017         Dai et al. 2019
Pohlen et al. 2017        Child et al. 2019

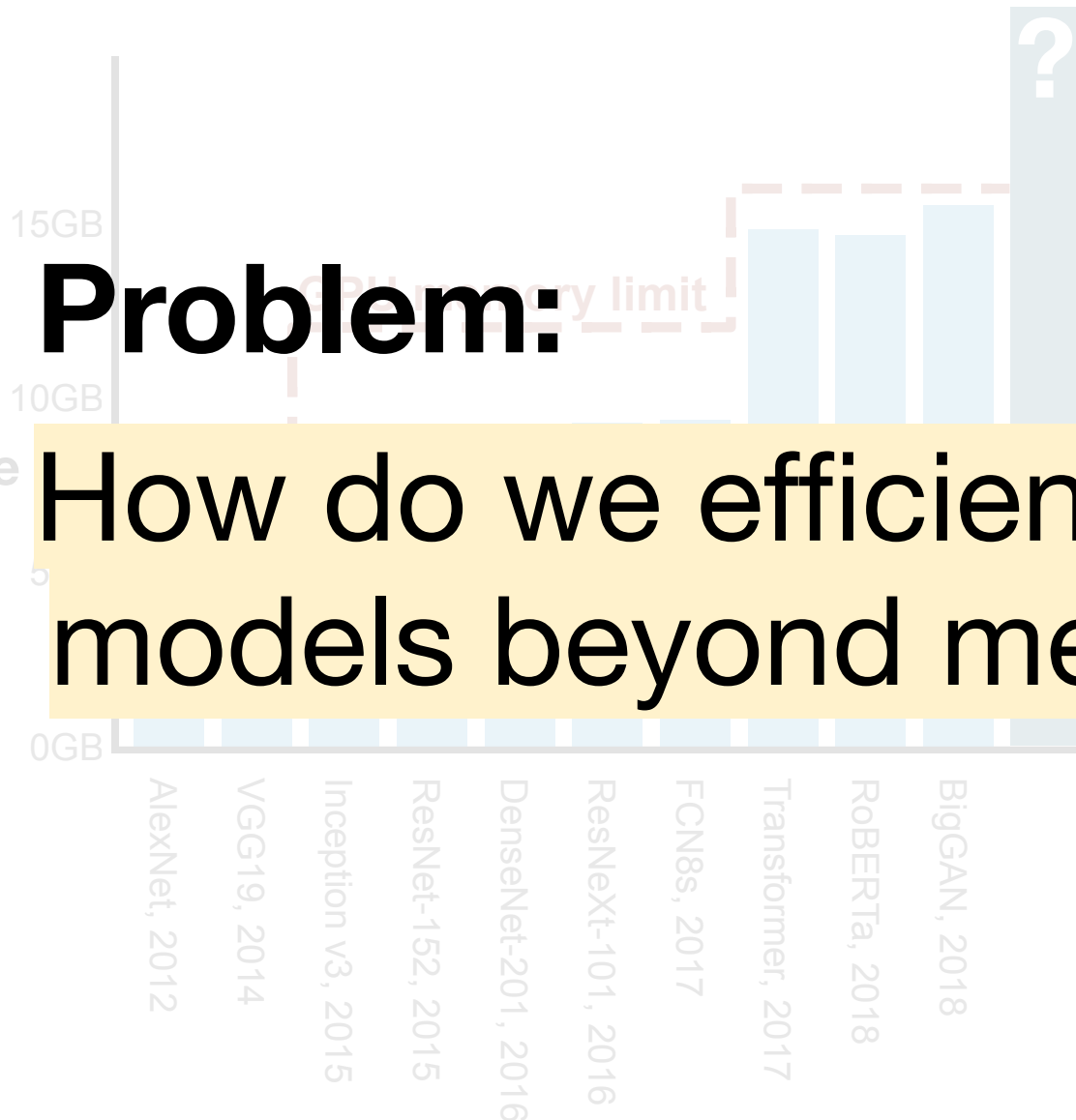**Limited GPU memory is slowing progress in new deep learning models!**

checkmateai.github.io

**Problem:**

State-of-the-art models have hit a **memory capacity wall**.

15GB

GPU memory limit

Cited memory as limiting factor

10GB

RAM usage per-GPU

How do we efficiently train large models beyond memory limits?

al. 2019
al. 2019
et al. 2019

0GB

AlexNet, 2012
VGG19, 2014
Inception v3, 2015
ResNet-152, 2015
DenseNet-201, 2016
ResNeXt-101, 2016
FCN8s, 2017
Transformer, 2017
RoBERTa, 2018
BigGAN, 2018

**Limited GPU memory is slowing progress in new deep learning models!**
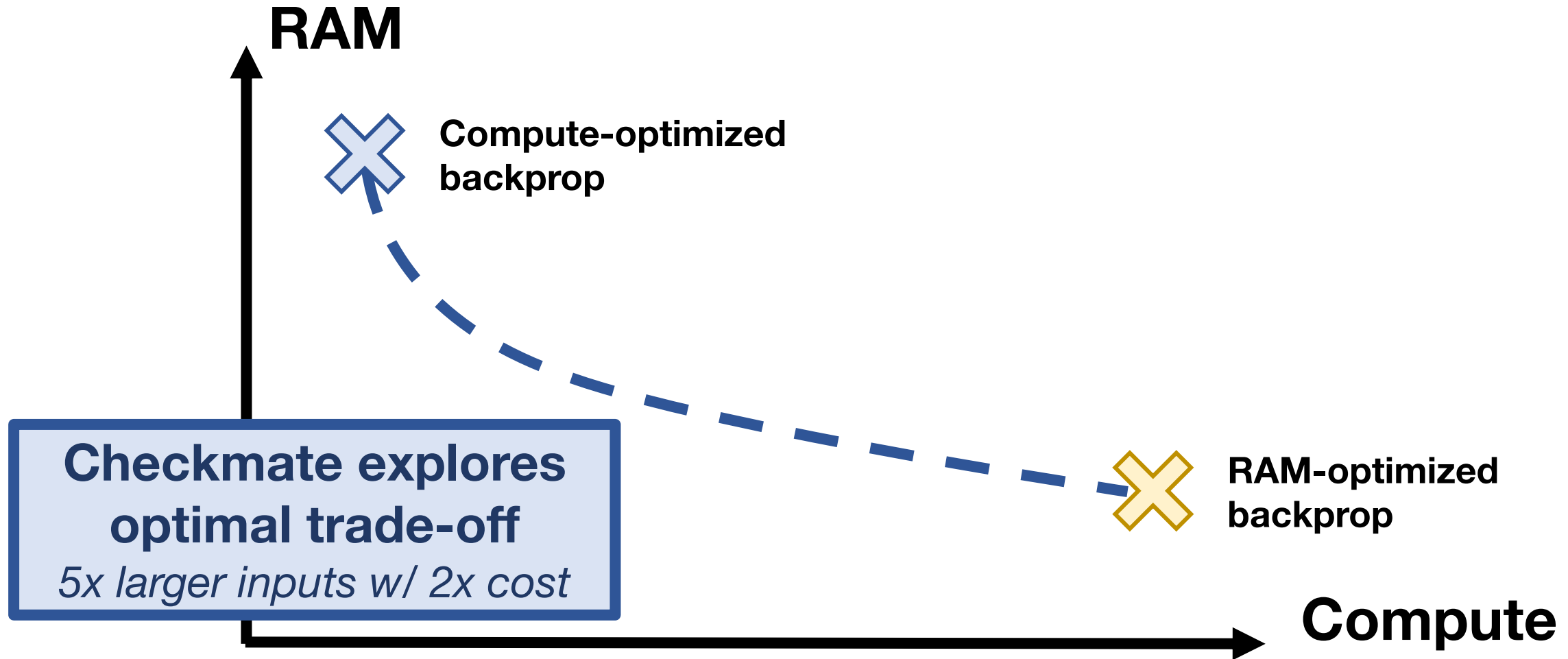
**Compute is outstripping DRAM capacity growth**

checkmateai.github.io

# Backprop is optimized for compute efficiency, not RAM usage

**RAM**

✕ **Compute-optimized backprop**

**Compute**

# **Ideal:** scalable algorithm for backprop that adapts to RAM constraints



**RAM**

**Compute-optimized backprop**

**RAM-optimized backprop**

**Compute**

checkmateai.github.io

# **This work:** optimal space-time tradeoff for backpropagation



**RAM**

✖ **Compute-optimized backprop**

**Checkmate explores optimal trade-off**
*5x larger inputs w/ 2x cost*
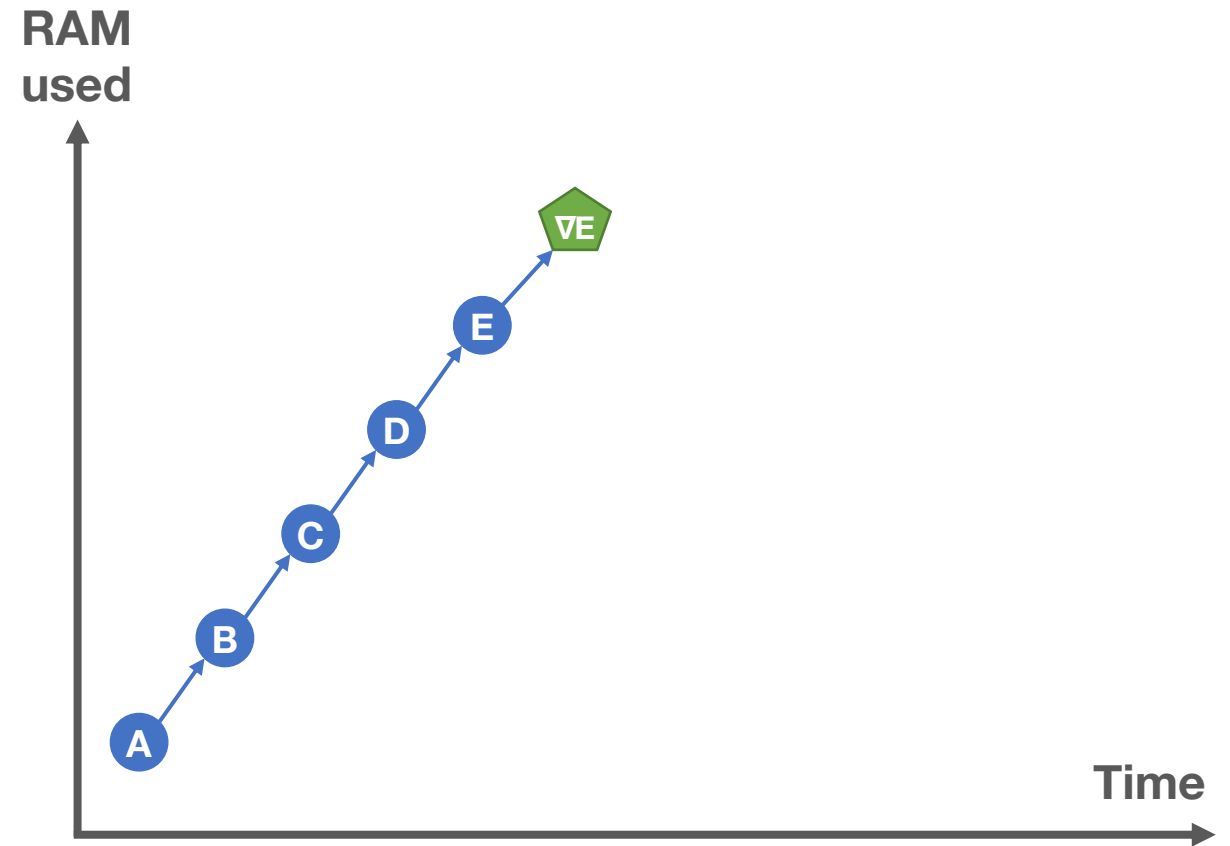
✖ **RAM-optimized backprop**

**Compute**

checkmateai.github.io

# RAM-hungry backprop policy
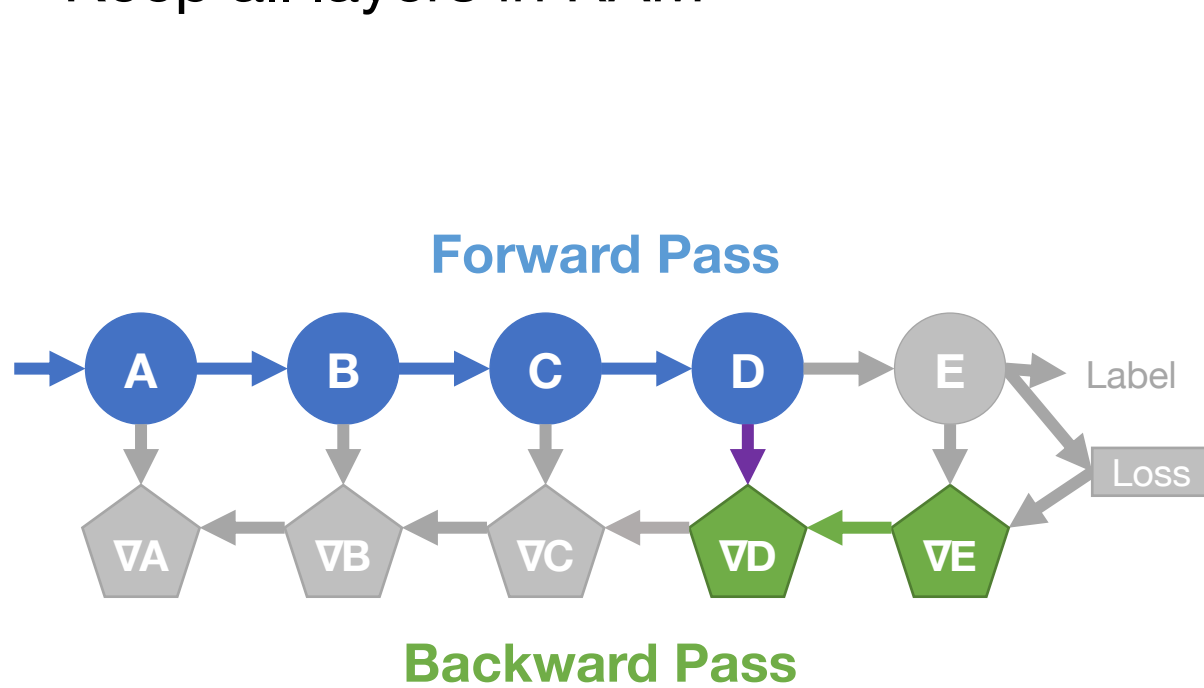
Keep all layers in RAM



RAM

Compute-optimized
backprop

Compute

# RAM-hungry backpropagation policy
## Keep all layers in RAM



Forward Pass

Backward Pass

RAM used

Time

# RAM-hungry backpropagation policy
## Keep all layers in RAM



Forward Pass

Backward Pass

RAM used
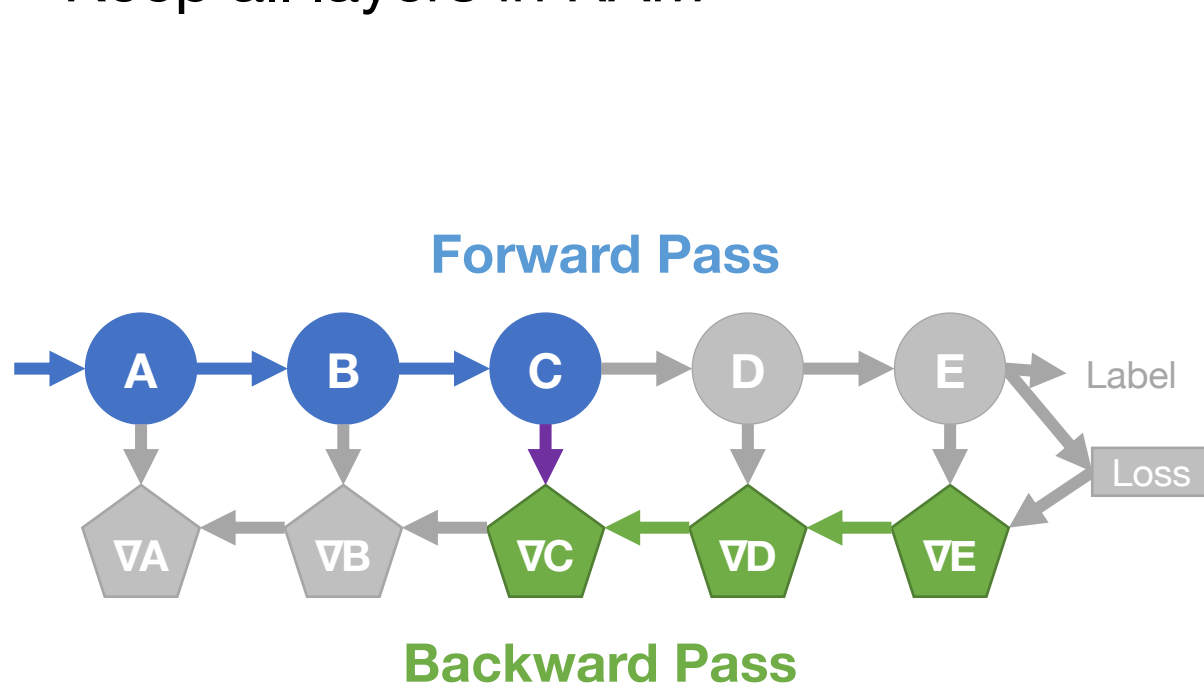
Time

# RAM-hungry backpropagation policy
Keep all layers in RAM

checkmateai.github.io
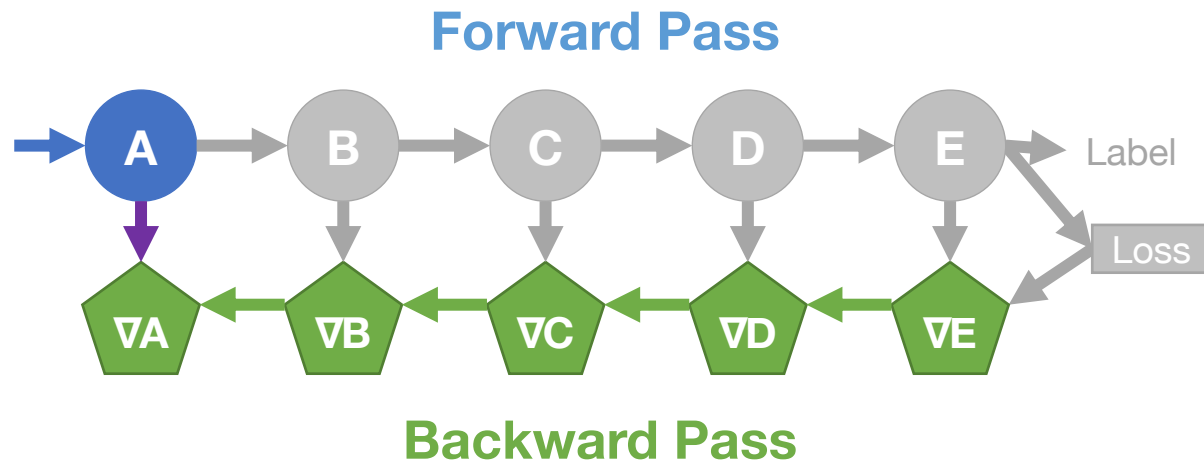
# RAM-hungry backpropagation policy
Keep all layers in RAM

checkmateai.github.io

# RAM-optimized backpropagation policy
Recompute all layers as needed

checkmateai.github.io

# RAM-optimized backpropagation policy

Recompute all layers

RAM used

**Peak RAM (no recomputation)**

**Forward Pass**

A → B → C → D → E → Label

Loss

∇A ← ∇B ← ∇C ← ∇D ← ∇E

**Backward Pass**

D
C
B
A

**Time**

# How can we use less memory?
Free early & recompute

checkmateai.github.io

# RAM-optimized backpropagation policy

Recompute <mark>all layers</mark>



**Forward Pass**

**Backward Pass**

**RAM used**
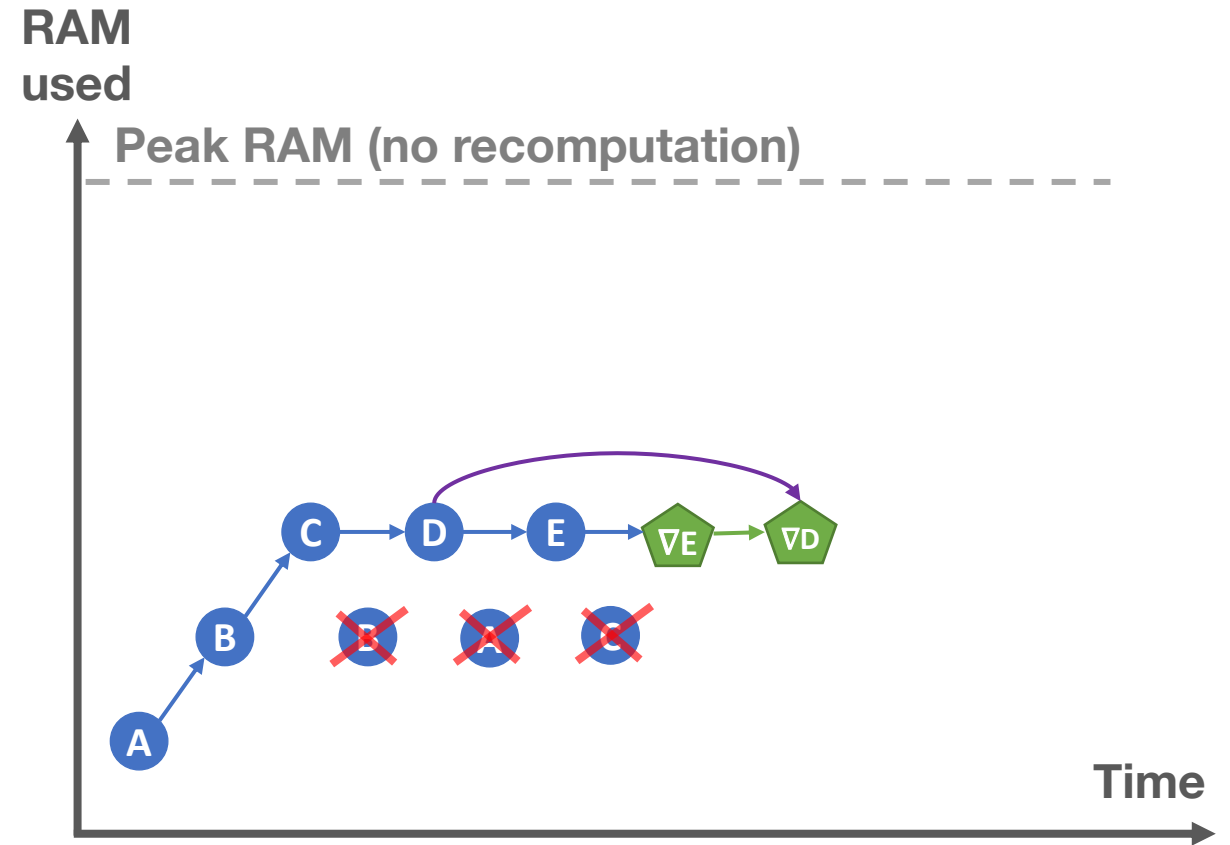
**Peak RAM (no recomputation)**

**Time**

# How can we use less memory?
## Free early & recompute

# RAM-optimized backpropagation policy

Recompute all layers



**How can we use less memory?**
Free early & recompute

checkmateai.github.io

# RAM-optimized backpropagation policy
Recompute <mark>all layers</mark>



**Forward Pass**

A → B → C → D → E → Label

Loss

∇A ← ∇B ← ∇C ← ∇D ← ∇E

**Backward Pass**

RAM used

**Peak RAM (no recomputation)**

**Peak RAM**

Time

# How can we use less memory?
Free early & recompute

checkmateai.github.io

# How to choose which layers to recompute?



**Forward Pass**

A → B → C → C → C → C → C → D → E → Loss

**Backward Pass**

∇A ← ∇B ← ∇C ← ∇C ← ∇C ← ∇C ← ∇C ← ∇D ← ∇E

checkmateai.github.io

# How to choose which layers to recompute?

**Forward Pass**



Label

checkmateai.github.io

# How to choose which layers to recompute?



**Compute:** $O(n)$ additional overhead

**RAM:** $O(\sqrt{n})$ RAM usage

checkmateai.github.io

**Challenges of heuristics:**

1. **Variable runtime per layer**

$10^6\times$
slower

Label

Label

$10^3\times$ more RAM

**Challenges of heuristics:**

1. Variable runtime per layer

2. Variable RAM usage per layer

checkmateai.github.io

**Challenges of heuristics:**

1. Variable runtime per layer

2. Variable RAM usage per layer

3. Real DNNs are non-linear

checkmateai.github.io

# Prior work is suboptimal in general setting!

**Greedy heuristic**

[Chen 2016]

[XLA authors 2017, 2020]

**Divide-and-conquer heuristic**

[Griewank 2000]

[Kowarz 2006]

[Siskind 2018]

[Kumar 2019]

**Optimal for specific architecture**

[Gruslys 2016]

[Feng 2018]

[Beaumont 2019]

*Challenges:*

1. **Variable runtime per layer**

2. **Variable RAM usage per layer**

3. **Real DNNs are non-linear**

checkmateai.github.io

# Can we ==optimally== trade-off RAM for compute?

*Let's be:*

**1. Hardware-aware**

**2. RAM-aware**

**3. DAG flexibility**

**RAM**

✖ Checkpoint every node

✖ Recompute all layers

**Compute**

checkmateai.github.io

# ✅ Checkmate

## A system for ==**optimal**== tensor rematerialization

**Hardware + RAM aware**

**Solve for 10s-1hr Train for 1mo**

**GPU, CPU, TPU support**

**Accurate cost model**

**Flexible search space**

**Optimal solver**
Integer Linear Program

**Near-optimal solver**
Two phase rounding

**Graph rewrite**

# ☑🤖 Checkmate

## A system for ==**optimal**== tensor rematerialization

```
┌─────────────┐          ┌───────────────────────────┐     ┌──────────┐
│  Accurate   │  ──────▶ │  ┌─────────────────────┐  │     │  Graph   │
│ cost model  │          │  │   Optimal solver    │  │ ──▶ │ rewrite  │
│             │          │  │ Integer Linear Prog. │  │     │          │
└─────────────┘          │  └─────────────────────┘  │     │          │
                         │                           │     │          │
┌─────────────┐          │  ┌─────────────────────┐  │     │          │
│  Flexible   │  ──────▶ │  │  Near-optimal solver │  │     │          │
│search space │          │  │ Two phase rounding   │  │     │          │
└─────────────┘          │  └─────────────────────┘  │     └──────────┘
                         └───────────────────────────┘
```
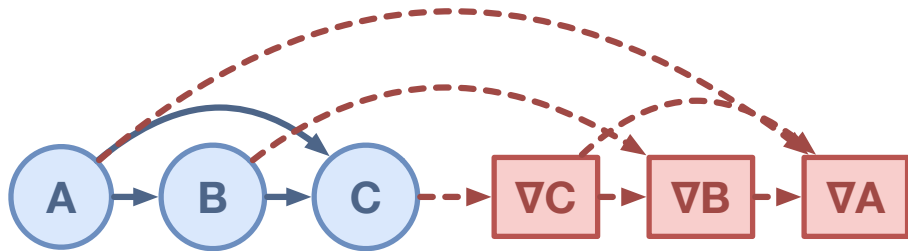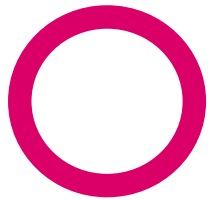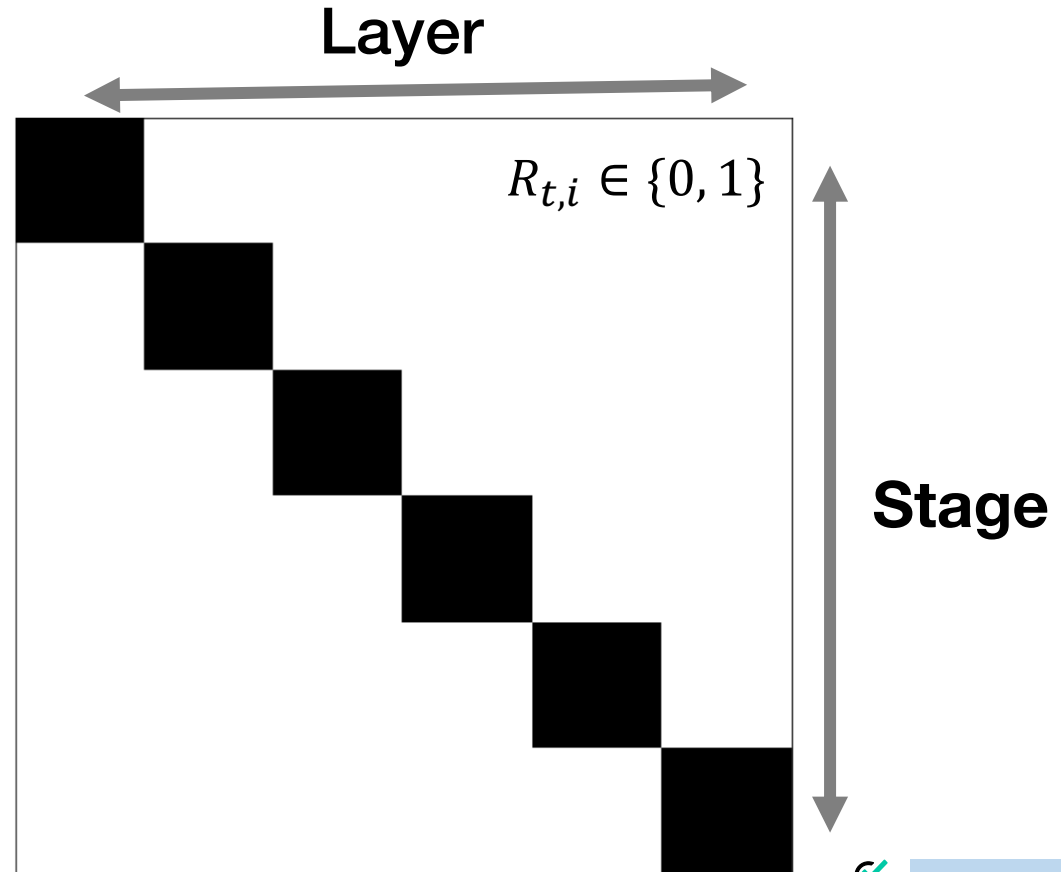
# ✅ Checkmate

## A system for ==optimal== tensor rematerialization

# ✅ Checkmate

## A system for ==**optimal**== tensor rematerialization



t=1

Layer

$R_{t,i} \in \{0, 1\}$

Stage

# ✅ Checkmate

## A system for ==**optimal**== tensor rematerialization

$R_{t,i} \in \{0, 1\}$

Layer

Stage

t=1  A  B  C  ∇C  ∇B  ∇A
t=2  A  B  C  ∇C  ∇B  ∇A
t=3  A  B  C  ∇C  ∇B  ∇A
t=4  A  B  C  ∇C  ∇B  ∇A
t=5  A  B  C  ∇C  ∇B  ∇A
t=6  A  B  C  ∇C  ∇B  ∇A

checkmateai.github.io

# ✅🤖 Checkmate

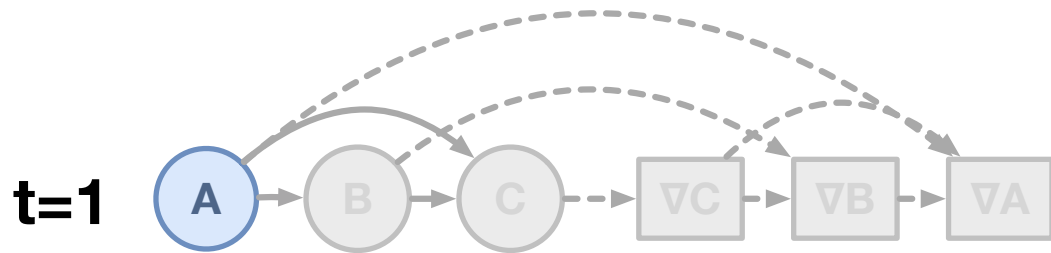## A system for <mark>**optimal**</mark> tensor rematerialization



Layer

Layer

$S_{t,i} \in \{0, 1\}$

$R_{t,i} \in \{0, 1\}$

t=1

t=2

t=3

t=4

t=5

t=6

Stage

# ✅🤖 Checkmate

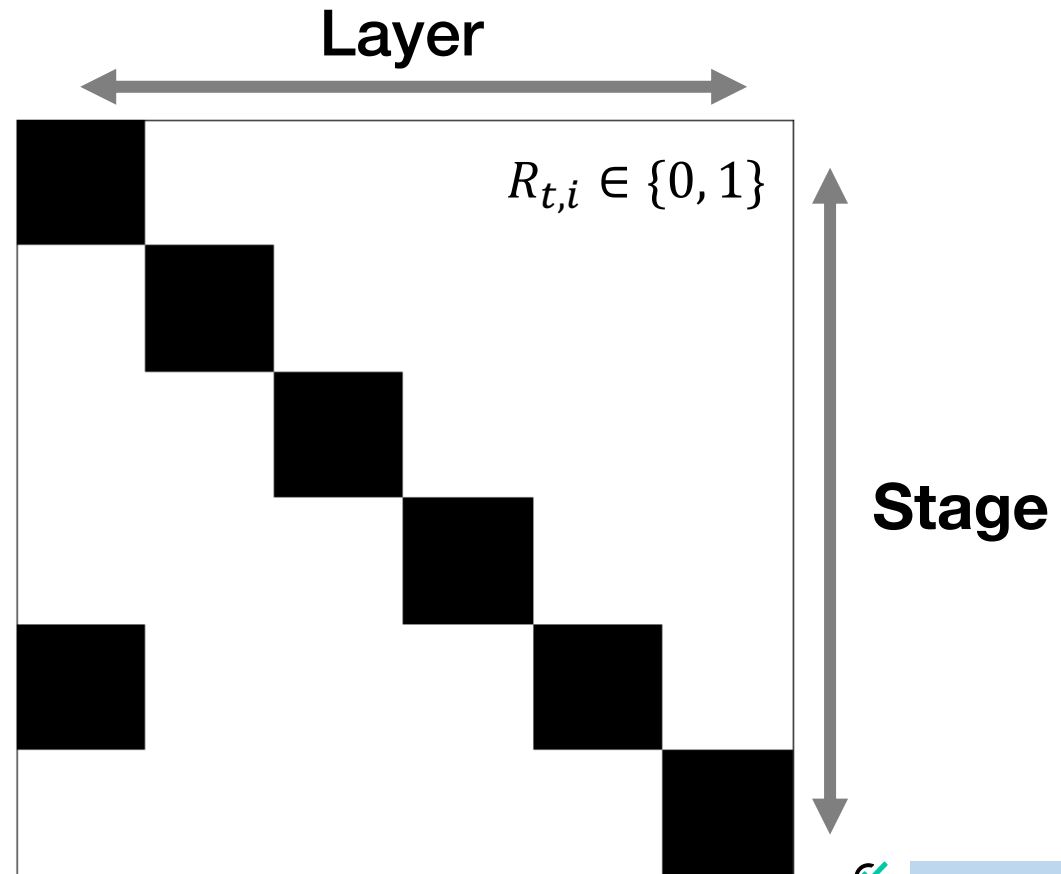## A system for **optimal** tensor rematerialization

Layer

Layer

$S_{t,i} \in \{0, 1\}$

$R_{t,i} \in \{0, 1\}$

t=1

t=2

t=3

t=4

t=5

t=6

Stage

**S = What is in memory?**

**R = What is computed?**

# ☑ Checkmate

## A system for ==<u>optimal</u>== tensor rematerialization



Accurate cost model

Flexible search space

**Optimal solver**
Integer Linear Program

Near-optimal solver
Two phase rounding

Graph rewrite

# 🤖 Checkmate

## A system for ==optimal== tensor rematerialization

**Minimize forward + backward cost**

$$\min_{S,R,U} \sum \sum C_i \, R_{t,i}$$

**Use R matrix to create linear objective**

**Decision variables**

$S_{t,i} \in \{0, 1\}$      Layer $i$ <u>s</u>tored for stage $t$

$R_{t,i} \in \{0, 1\}$      Layer $i$ (<u>re</u>)computed in stage $t$

# ✓🤖 Checkmate

## A system for ==optimal== tensor rematerialization

Minimize forward + backward cost

$$\min_{S,R,U} \sum \sum C_i \, R_{t,i}$$

**Decision variables**

$S_{t,i} \in \{0, 1\}$ — Layer $i$ <u>s</u>tored for stage $t$

$R_{t,i} \in \{0, 1\}$ — Layer $i$ (<u>re</u>)computed in stage $t$

**Correctness**

$$R_{t,j} \leq R_{t,i} + S_{t,i}$$

$$S_{t,i} \leq R_{t-1,i} + S_{t-1,i}$$

"A layer's dependencies must be computed before evaluation"

"A layer must be computed before it can be stored in RAM"

# ✅🤖 Checkmate

## A system for ==optimal== tensor rematerialization

**Minimize forward + backward cost**

$$\min_{S,R,U} \sum \sum C_i \, R_{t,i}$$

**Decision variables**

| | |
|---|---|
| $S_{t,i} \in \{0,1\}$ | Layer $i$ <u>s</u>tored for stage $t$ |
| $R_{t,i} \in \{0,1\}$ | Layer $i$ (<u>re</u>)computed in stage $t$ |
| $U_{t,i} \in \mathbb{R}_+$ | Memory <u>u</u>sage in stage $t$ |

**Correctness**

$$R_{t,j} \le R_{t,i} + S_{t,i}$$

$$S_{t,i} \le R_{t-1,i} + S_{t-1,i}$$

**Memory limit**

$$U_{t,k} \le \text{budget}, \dots$$

**Constrain memory via an implicit variable to model memory usage at each stage**

✅🤖 checkmateai.github.io

# ✅ Checkmate

## A system for <mark>optimal</mark> tens

Constrain memory at all execution steps

$$U_{t,k} \leq \text{budget}$$

More checkpoints → more initial memory

$$U_{t,0} = \sum M_i S_{t,i}$$

$$U_{t,k+1} = U_{t,k} - \boxed{\text{GC Deps}[v_k]} + \boxed{M_{k+1} R_{t,k+1}}$$

**Minimize forward + backward cost**

$$\min_{S,R,U} \sum \sum C_i\, R_{t,i}$$

**Correctness**

$$R_{t,j} \leq R_{t,i} + S_{t,i}$$

$$S_{t,i} \leq R_{t-1,i} + S_{t-1,i}$$

**Memory limit**

$$U_{t,k} \leq \text{budget}, \dots$$

**Memory accounting details in paper**
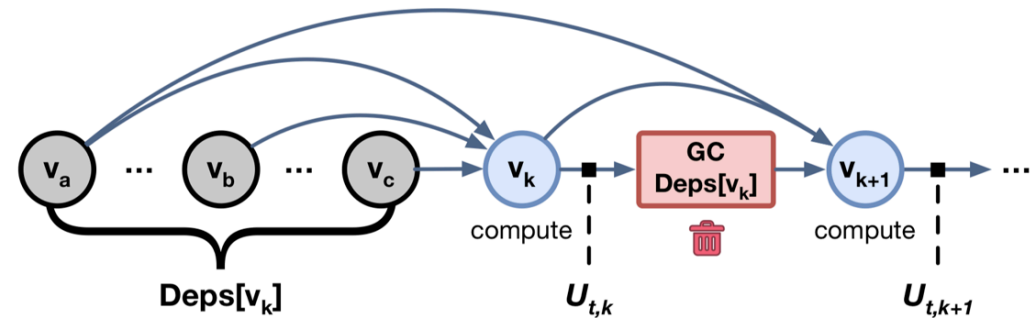
Constrain me
to model memory usage at each stage



Memory Usage

# ✓ Checkmate

## A system for ==optimal== tensor rematerialization

**Minimize forward + backward cost**

$$\min_{S,R,U} \sum \sum C_i \, R_{t,i}$$

**Correctness**

$$R_{t,j} \leq R_{t,i} + S_{t,i}$$

$$S_{t,i} \leq R_{t-1,i} + S_{t-1,i}$$

**Memory limit**

$$U_{t,k} \leq \text{budget}, \dots$$

**How long is the solve time?**

9 hours 😳

# ☑🤖 Checkmate

## A system for ==**optimal**== tensor rematerialization

**Minimize forward + backward cost**

$$\min_{S,R,U} \sum \sum C_i \, R_{t,i}$$

**Correctness**

$$R_{t,j} \leq R_{t,i} + S_{t,i}$$

$$S_{t,i} \leq R_{t-1,i} + S_{t-1,i}$$

**Memory limit**

$$U_{t,k} \leq \text{budget}, \ldots$$

**Partition schedule into frontier-advancing stages**
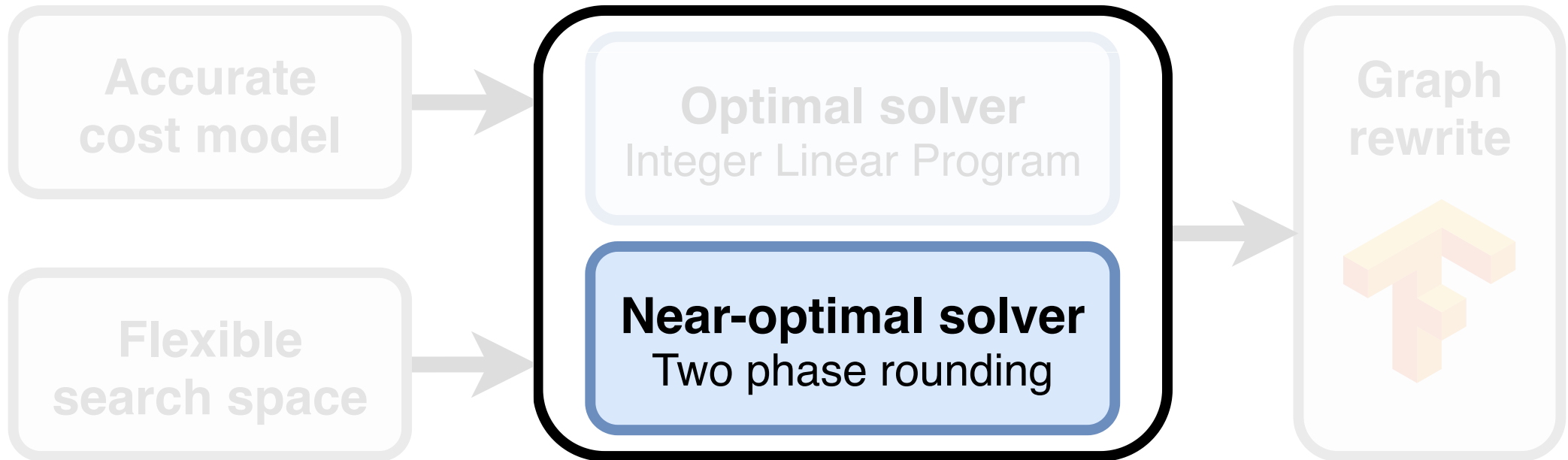
9 hours → 0.2 seconds

**Tractability**

$$R_{t,t} = 1$$

$$R, S, U \text{ lower triangular}$$

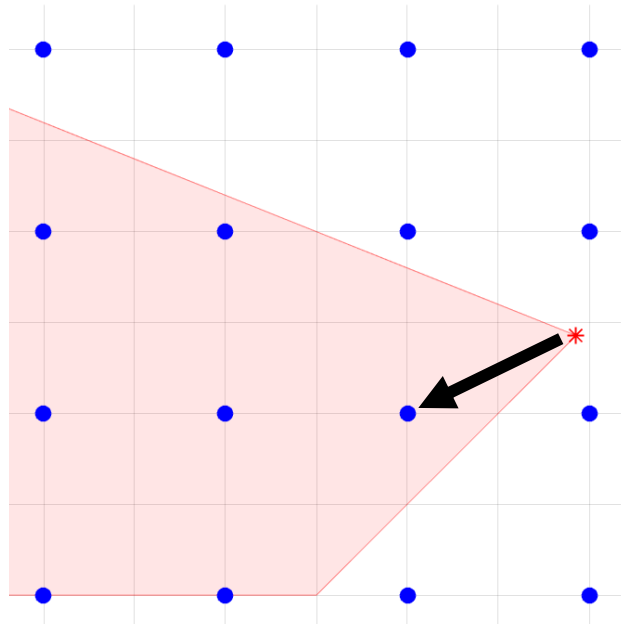Prunes n! permutations of nodes

# ☑🤖 Checkmate

## A system for ==optimal== tensor rematerialization

**Accurate cost model**

**Flexible search space**

**Optimal solver**
Integer Linear Program

**Near-optimal solver**
Two phase rounding

**Graph rewrite**

# ILP optimization is NP-hard (combinatorial search)

Polynomial-time approximation?



1. Relax boolean constraints
2. Solve LP
3. Round solution

**How to maintain feasibility?**

**Insight: Given S, optimal R easy to compute**
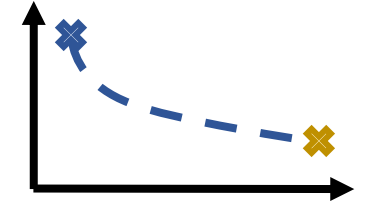
**Proposed method: Two-Phase Rounding**
Round S, solve other variables optimally
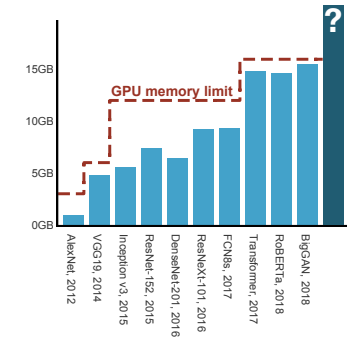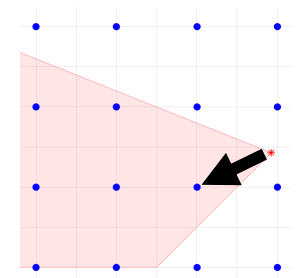
# **Evaluation:** Questions

1. What is the memory vs compute trade-off?



2. How much can we increase batch/model size?



3. How well does two-phase rounding do?

# Evaluation: What is the memory vs compute trade-off?



U-Net, batch size 32

MobileNet, batch size 512

Compute Overhead (x)

Best heuristic

Checkmate

1.2x speedup!

GPU Memory available (GB)

# Evaluation: How much can we increase batch size?

Layer

Stage

$R_{t,i}$

**No rematerialization**
**Batch size 167**

Layer

$R_{t,i}$

**Square root heuristic**
**Batch size 197**

==1.18x larger!==

Layer

$R_{t,i}$

**Checkmate**
**Batch size 289**

==1.73x larger!==  10 sec solve

checkmateai.github.io

# **Evaluation:** How much can we increase batch size?



5.1x larger!

1.73x larger!

**Normalized batch size**

- 5x
- 4x
- 3x
- 2x
- 1x
- 0x

**U-Net** | **FCN8** | **SegNet** | **VGG19** | **ResNet50** | **MobileNet**

U-Net: 16, 18, 35, 61
FCN8: 29, 51, 51, 60
SegNet: 21, 33, 43, 62
VGG19: 167, 197, 266, 289
ResNet50: 98, 116, 199, 225
MobileNet: 215, 452, 640, 1105

■ Checkpoint all ■ AP √n ■ Lin. greedy ■ Checkmate (ours)

checkmateai.github.io

# Evaluation: How much can we increase batch size?



5.1x larger!

1105

5x

*Ongoing work: BERT
2.3x larger batch size over TF2.0
Train BERT-Large w/o model parallelism

2x

16  18    29  5  5    21  33    167  197  26  2    98  116    215

1x

0x

U-Net    FCN8    SegNet    VGG19    ResNet50 MobileNet

■ Checkpoint all   ■ AP √n   ■ Lin. greedy   ■ Checkmate (ours)

checkmateai.github.io

# **Evaluation:** How well does 2P rounding approximate ILP?

| | AP $\sqrt{n}$ | AP greedy | Griewank $\log n$ | Two-phase LP rounding |
|---|---|---|---|---|
| MobileNet | $1.14\times$ | $1.07\times$ | $7.07\times$ | $\mathbf{1.06\times}$ |
| VGG16 | $1.28\times$ | $1.06\times$ | $1.44\times$ | $\mathbf{1.01\times}$ |
| VGG19 | $1.54\times$ | $1.39\times$ | $1.75\times$ | $\mathbf{1.00\times}$ |
| U-Net | $1.27\times$ | $1.23\times$ | - | $\mathbf{1.03\times}$ |
| ResNet50 | $1.20\times$ | $1.25\times$ | - | $\mathbf{1.05\times}$ |

**Within 6% of optimal cost** (geomean)

43x speedup for ResNet50
440x speedup for MobileNet

# Checkmate

**Code and paper:**
checkmateai.github.io

**Email me:**
parasj@berkeley.edu

**Key ideas:**

- GPU memory limits are preventing the development of new deep learning models.

- We present the first general solution for optimal & near-optimal graph rematerialization.

- Formulation supports **arbitrary DAGs** and is both **hardware-aware** and **memory-aware**

- Integration with just **one line of code**

```
train_iteration = checkmate.compile_tf2(
    model, loss, optimizer, input_shape, label_shape)
```

Paras Jain, Ajay Jain, Ani Nrusimha, Amir Gholami, Pieter Abbeel, Kurt Keutzer, Ion Stoica, Joseph Gonzalez

Berkeley
UNIVERSITY OF CALIFORNIA