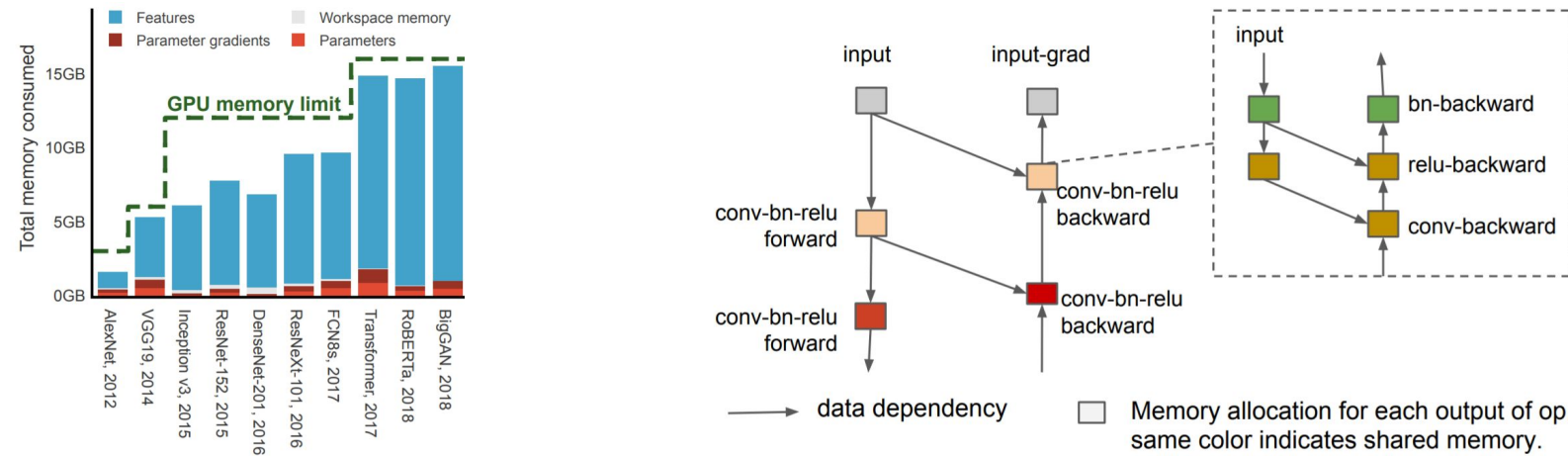


# Dynamic Tensor Rematerialization

Marisa Kirisame\*, Steven Lyubomirsky\*, Altan Haan\*, Jennifer Brennan, Mike He, Jared Roesch, Tianqi Chen, Zachary Tatlock

## Checkpointing: Trade Time for Space



“The memory wall” (Jain et al., 2020) Past approaches: Static plan (Chen et al., 2016)

## DTR: Checkpointing is Caching

### PerformOp(op, args):

Note: Performs op(args), rematerializing any evicted arguments. Wraps every operator invocation.

Exclude members of args from eviction for any evicted arg in args:

Rematerialize(arg)

buf := AllocateBuffer(size(op(args)))

res := call op(args), store into buf

Permit eviction for members of args again

Update metadata for args and res

return res

### PerformEviction():

Free the tensor chosen by the heuristic

### Rematerialize(t):

op, args := operator and arguments that produced t (from metadata)  
return PerformOp(op, args)

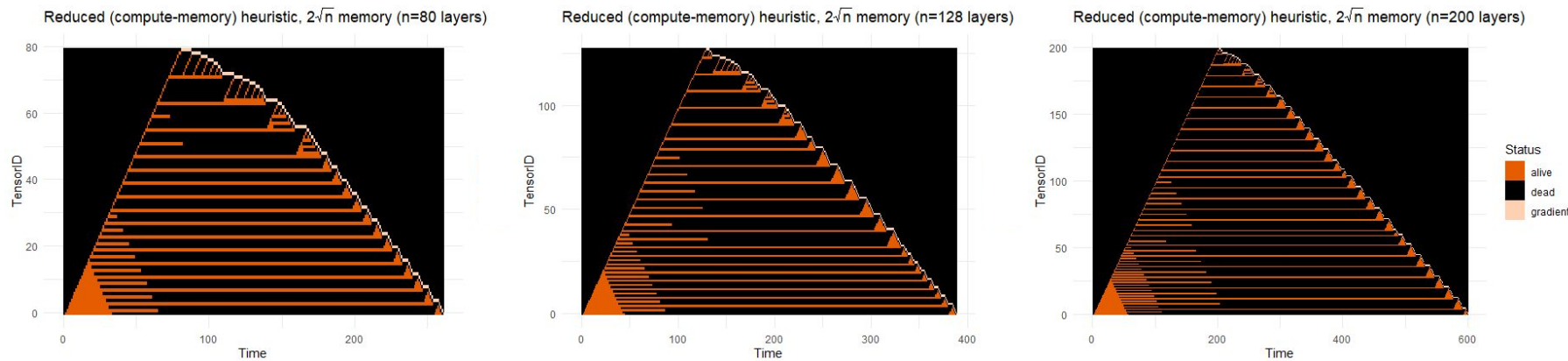
### AllocateBuffer(b):

Note: Wraps every memory allocation.  
while available memory < b:  
PerformEviction()  
return new buffer of size b

### Deallocate(t):

Note: Wraps every tensor deallocation.  
Heuristic decides policy for t (e.g., free permanently or simply evict)

## Theoretical Results



Train N-layer FF network in  $\Omega(\sqrt{N})$  memory and  $O(N)$  operations! No static planning!

## Execution Trace

Computing  $t_7$  with memory budget 4:

$t_7 = \text{PerformOp}(op_7, t_5, t_6)$

[ $t_5, t_6$  become unevictable]

Rematerialize( $t_5$ )

$t_5 = \text{PerformOp}(op_5, t_3)$

[ $t_3$  becomes unevictable]

AllocateBuffer( $t_5.size$ )

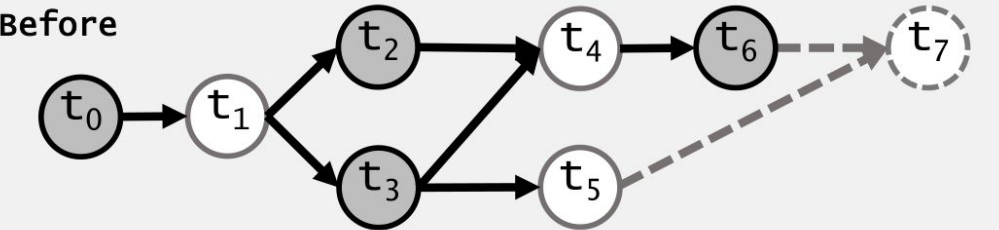
PerformEviction() #eg,  $t_2$

[ $t_3$  becomes evictable]

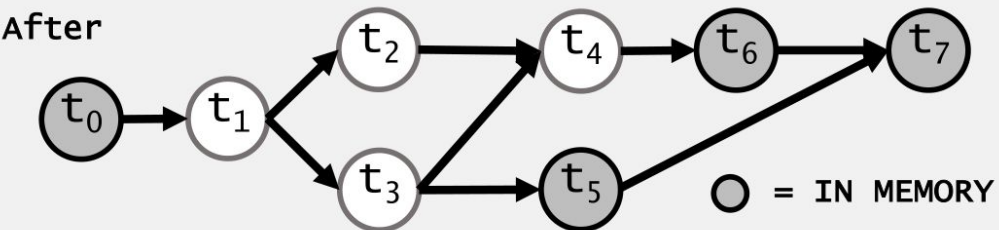
AllocateBuffer( $t_7.size$ )

PerformEviction() #eg,  $t_3$

Before



After



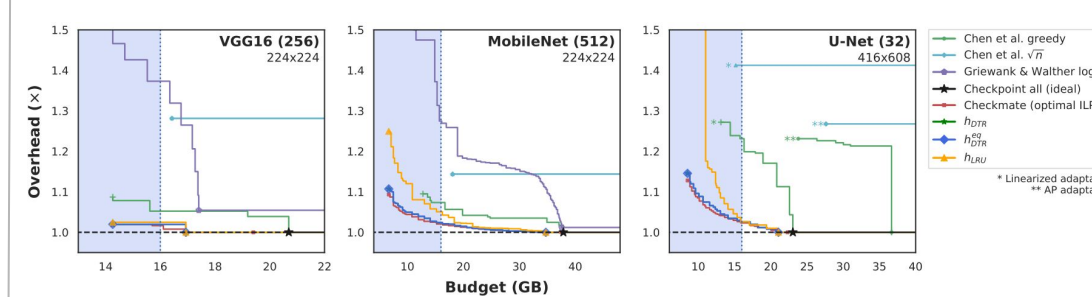
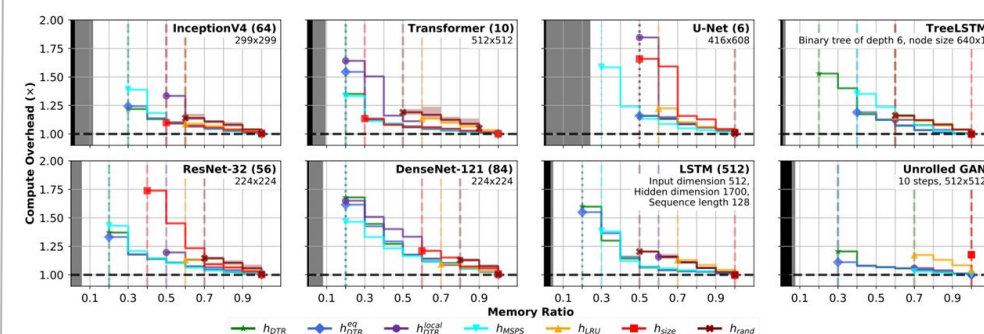
## Heuristics: The Brains of DTR

### Cost over dependencies

$$h_{\text{DTR}} \stackrel{\text{def}}{=} \frac{c_0(t) + \sum_{t' \in e^*(t)} c_0(t')}{m(t) \cdot s(t)}$$

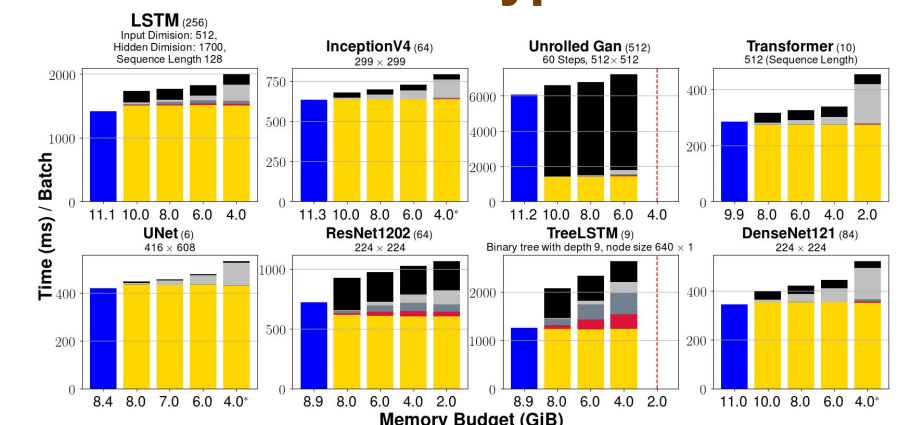
Size

Staleness



Near-optimal performance in simulation!

## Prototype



Few hundred lines in PT, limited overhead