# Faucet: a user-level, modular technique for flow control in dataflow engines

Andrea Lattuada
Systems Group,
ETH Zürich

Frank McSherry
Unaffiliated

Zaheer Chothia
Systems Group,
ETH Zürich

**Problem**        RAM exhaustion due to buffered
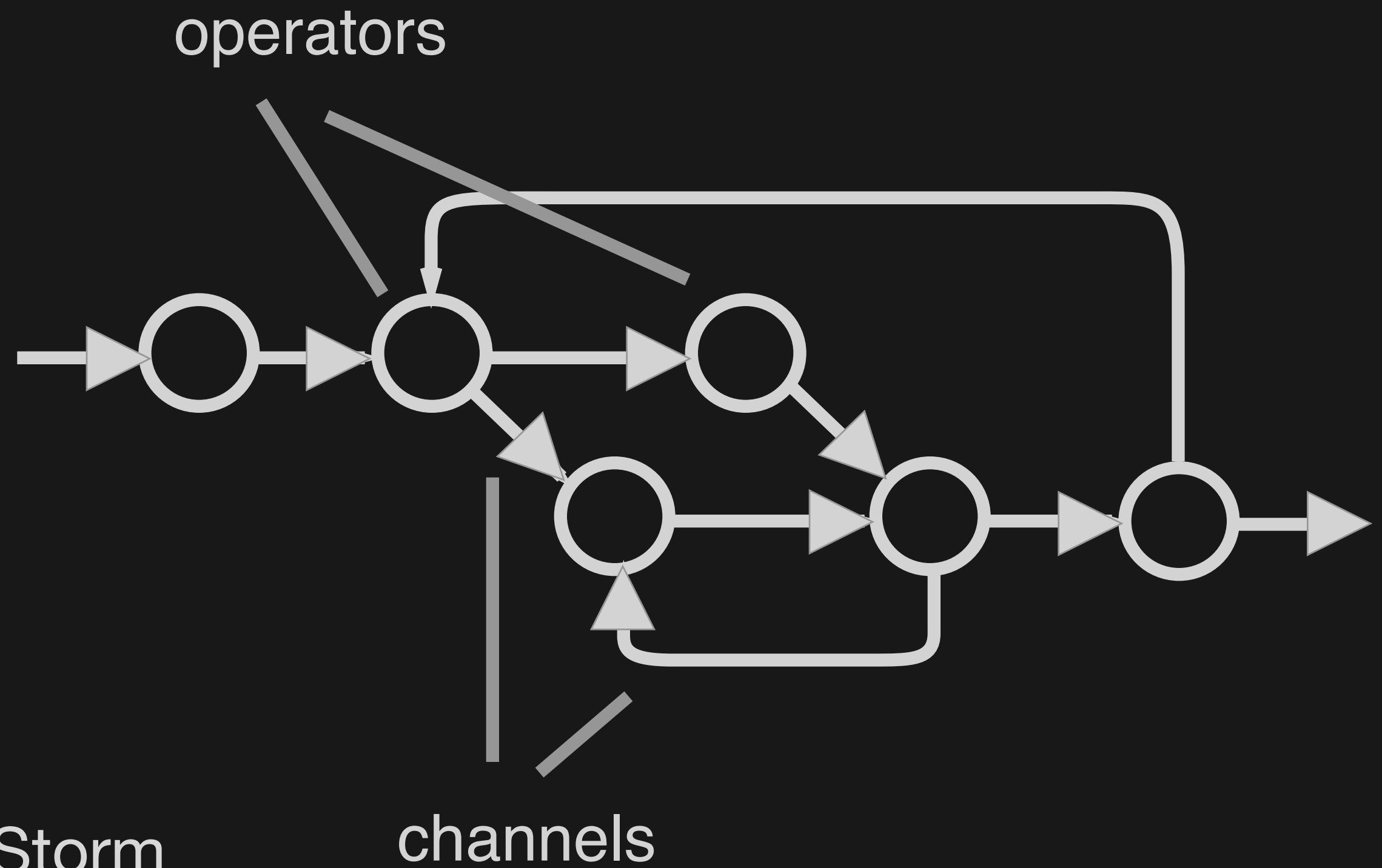*intermediate results*

**Our Solution**    • no system-level general strategy
• *application-driven* scheduling

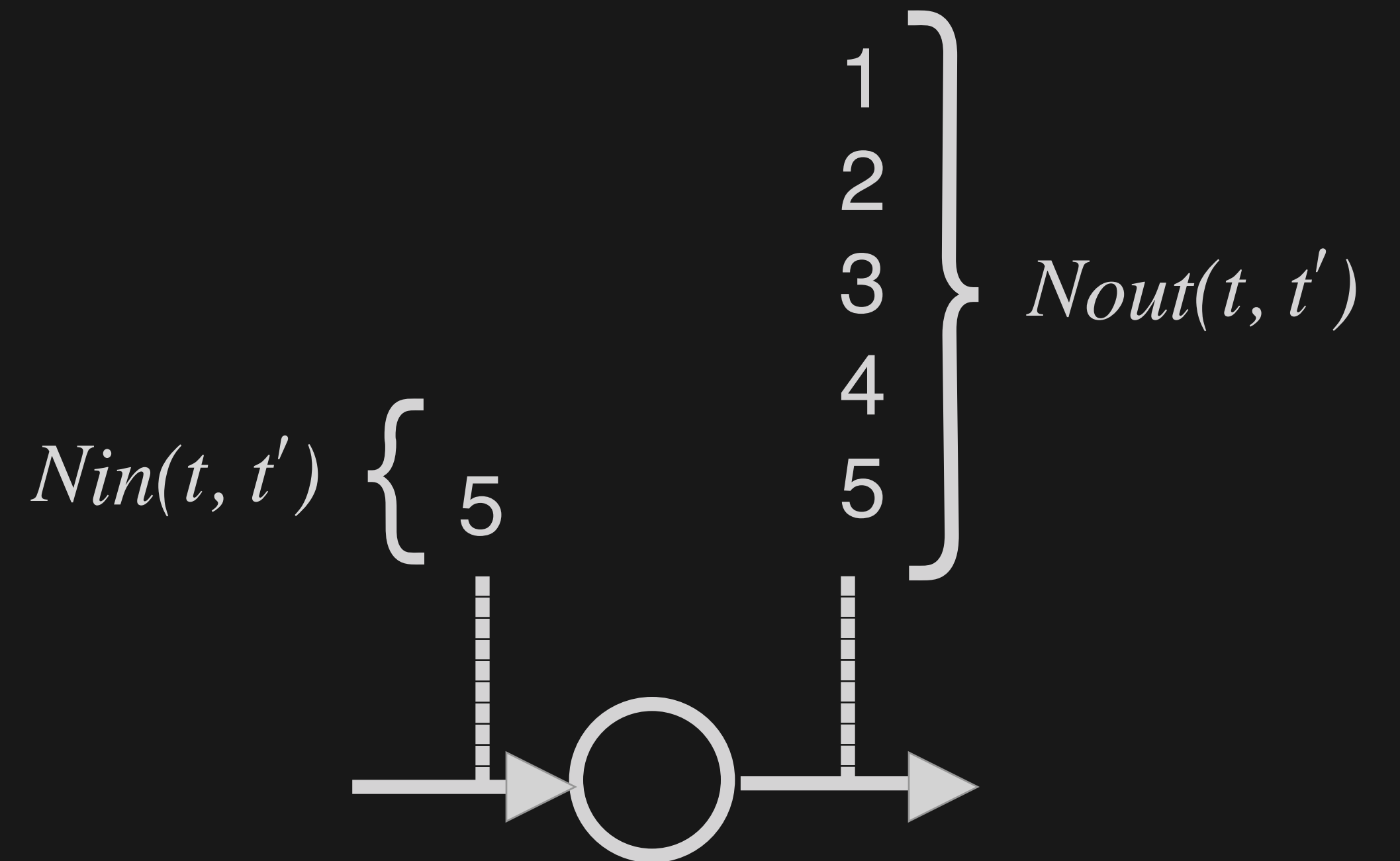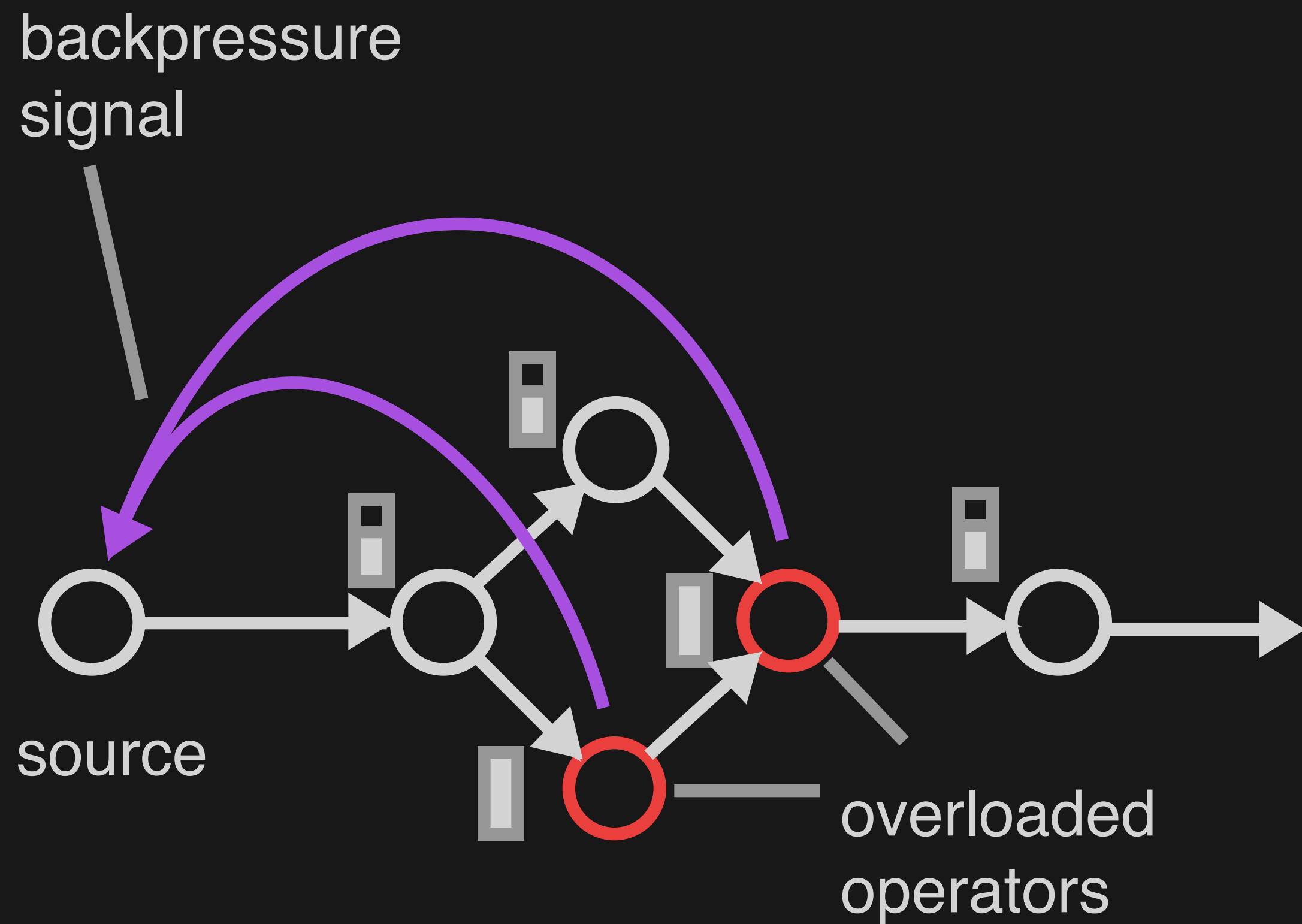10-100x memory savings for 15-25% runtime overhead

# Dataflow model



operators

channels

Storm

Flink

Naiad

## Rate imbalance
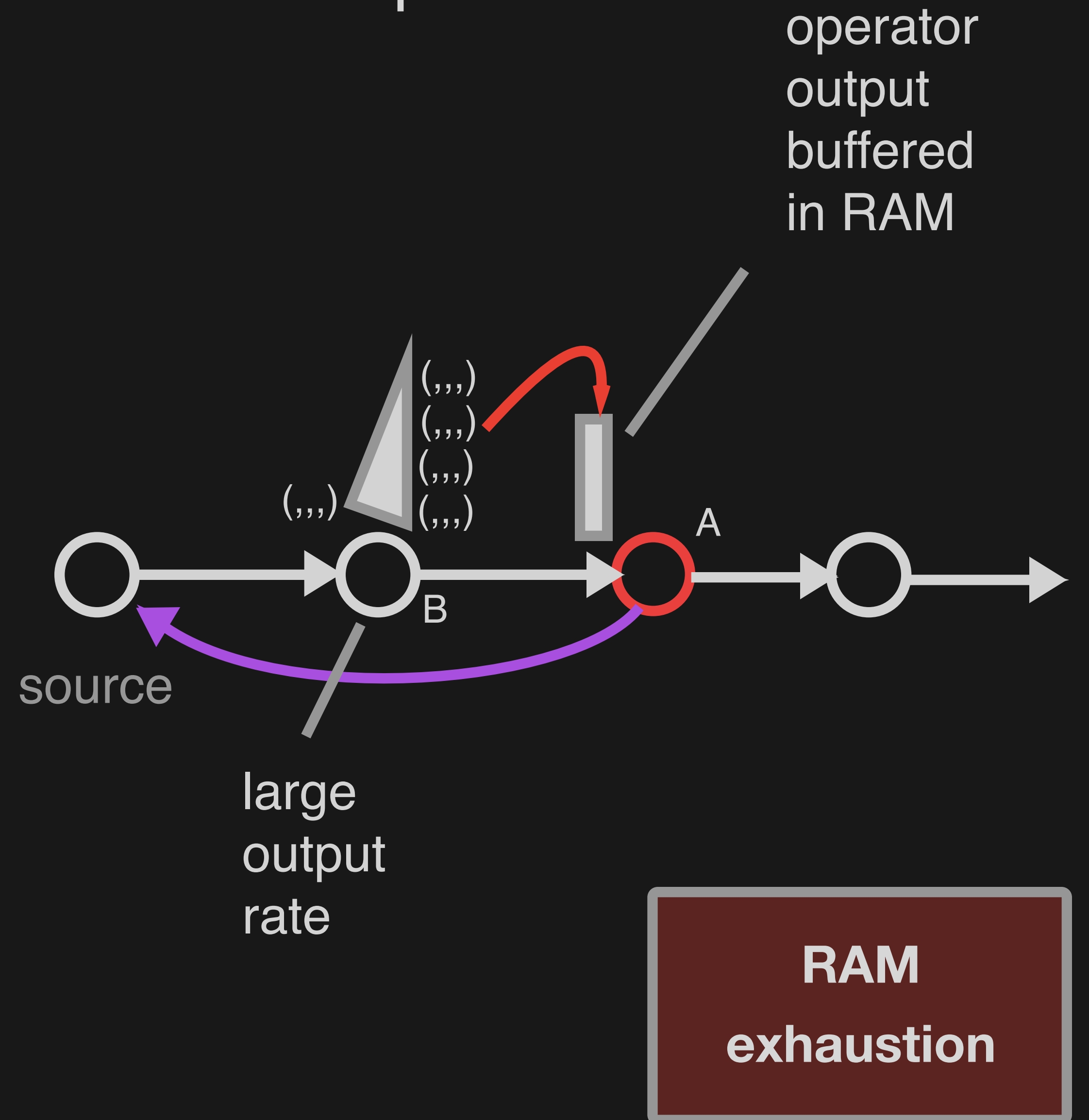


$$Nout(t, t')$$

$$Nin(t, t')$$

```
flat_map(|x| [1, …, x])
```

3

# Existing approach #1 - Source backpressure

backpressure signal

operator output buffered in RAM

source

overloaded operators

source

(,,,)
(,,,)
(,,,)
(,,,)
(,,,)

A

B

large output rate

Storm

Heron

Spark streaming

**RAM exhaustion**

# Existing approach #2 - Edge-by-edge backpressure

## similar to TCP flow control



backpressure
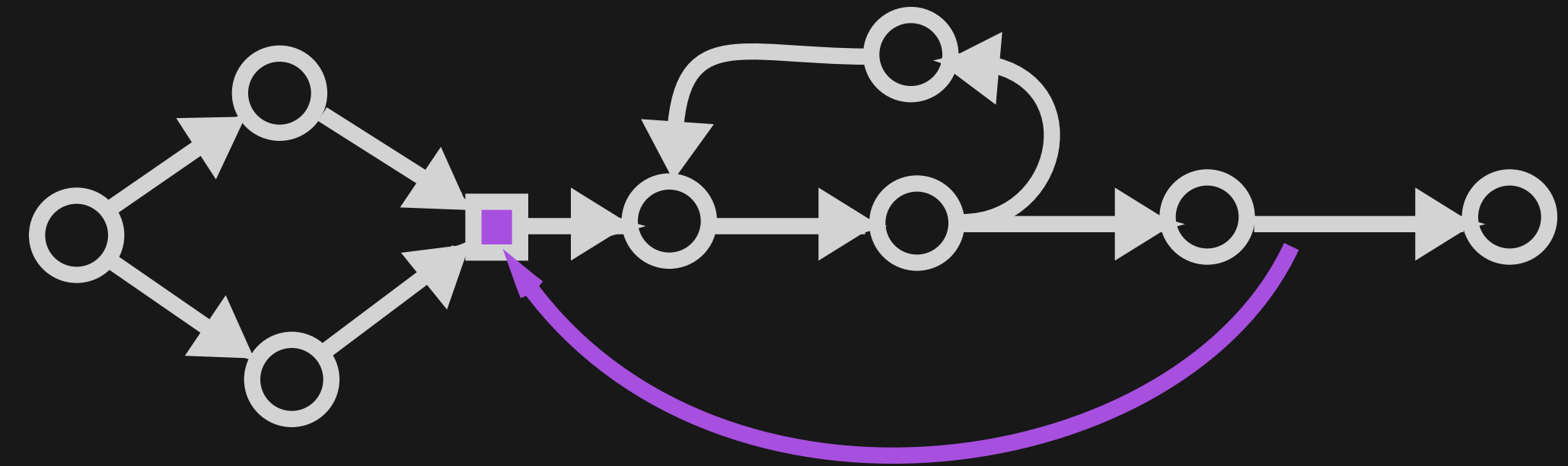signal

overloaded
operator

source

stopped by F

stopped by G

H

G

F

source

deadlock

Akka Streams

Flink

based on *Timely Dataflow's concepts*

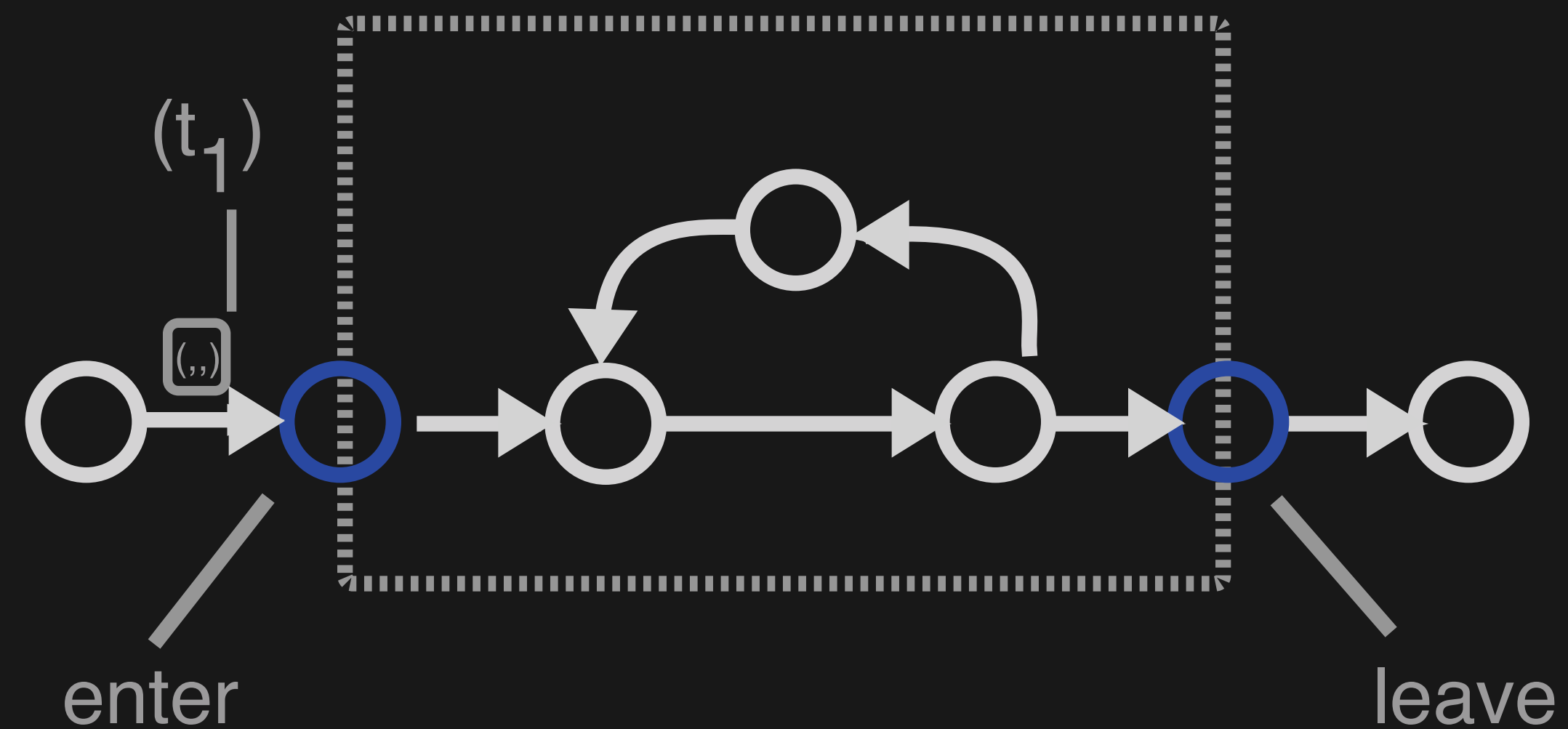- no fine-grained signal
- *track completion* of a batch
  of tuples

control scheduling to *limit intermediate results*

6

**Scopes**        nested operator structure

**Timestamps**    tuple metadata

**Scopes**      nested operator structure

**Timestamps**      tuple metadata

**Scopes**       nested operator structure
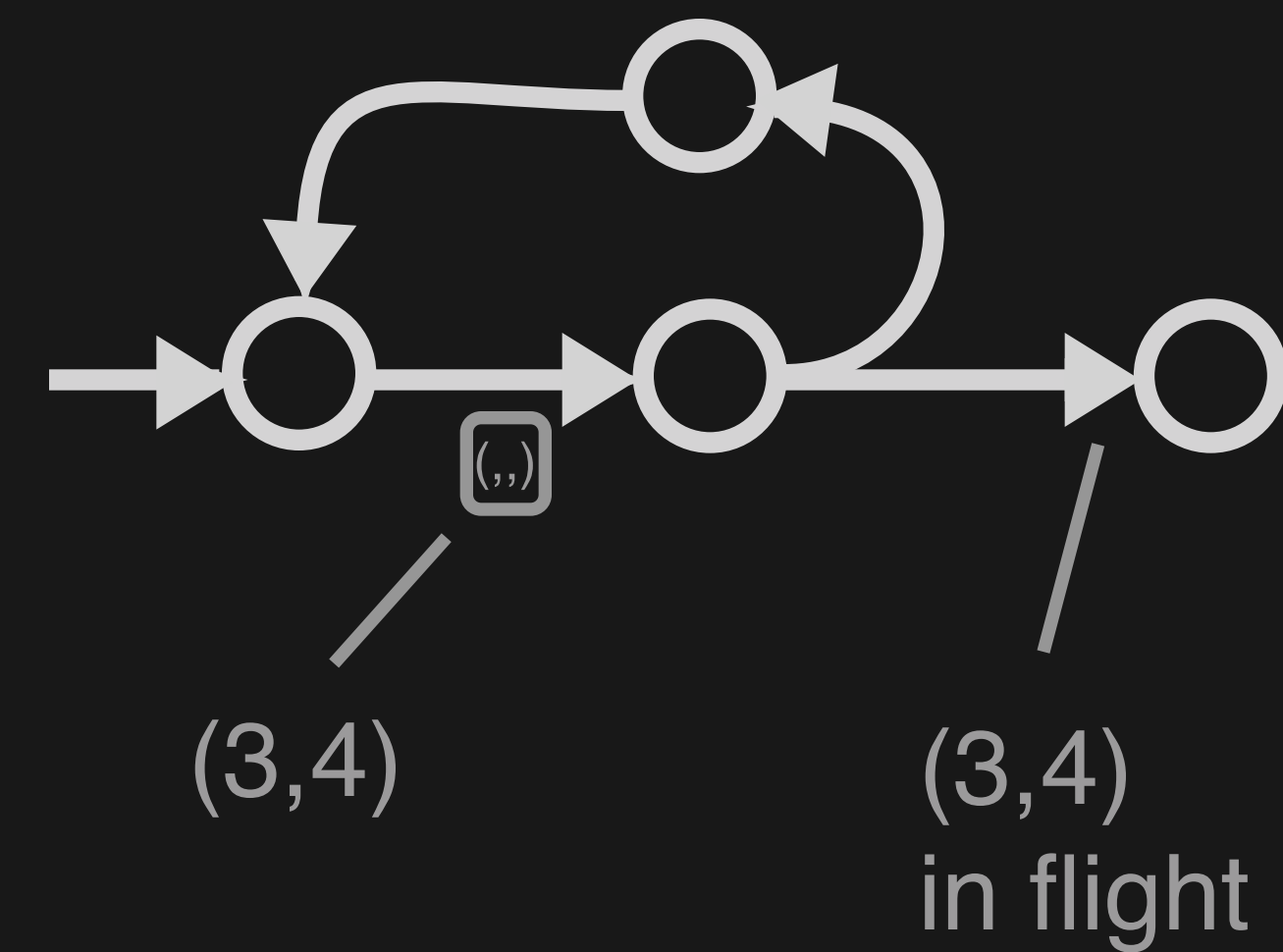
**Timestamps**   tuple metadata

**Scopes**        nested operator structure
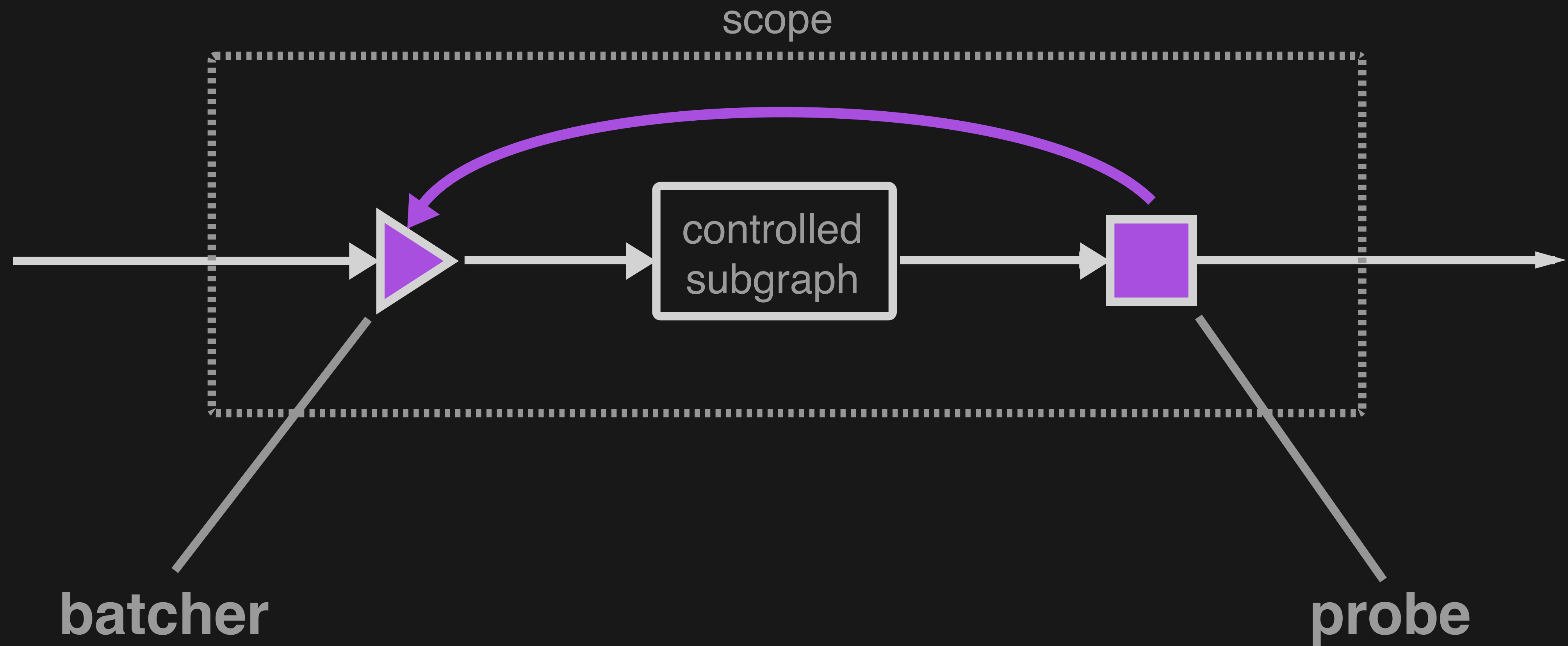
**Timestamps**    tuple metadata

**Progress Tracking**

tracks pending timestamps

# Faucet - Track batches of *intermediate results*



scope

controlled subgraph

**batcher**

**probe**

8

# Faucet - Track batches of *intermediate results*



scope

batcher

probe

controlled
subgraph

$(t_e)$

8

# Faucet - Track batches of *intermediate results*
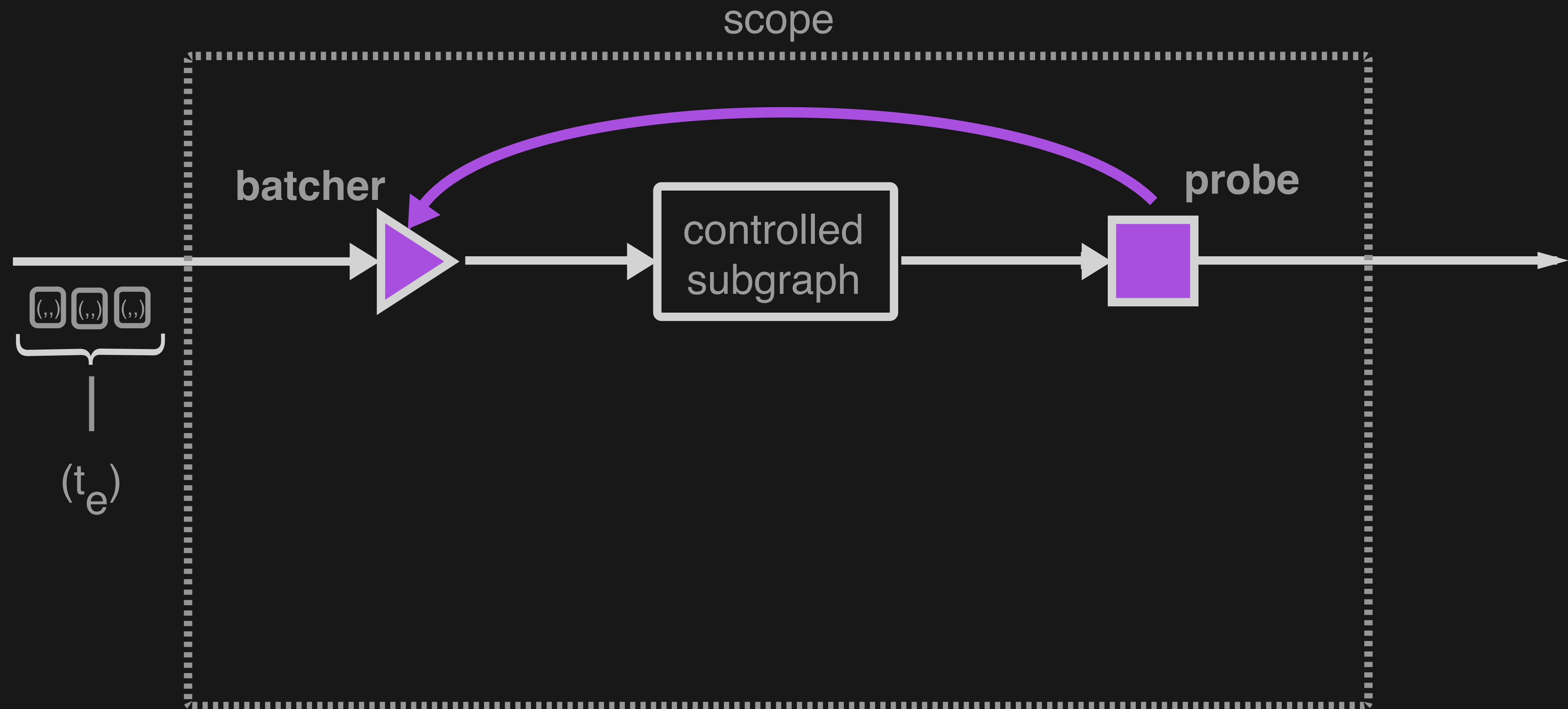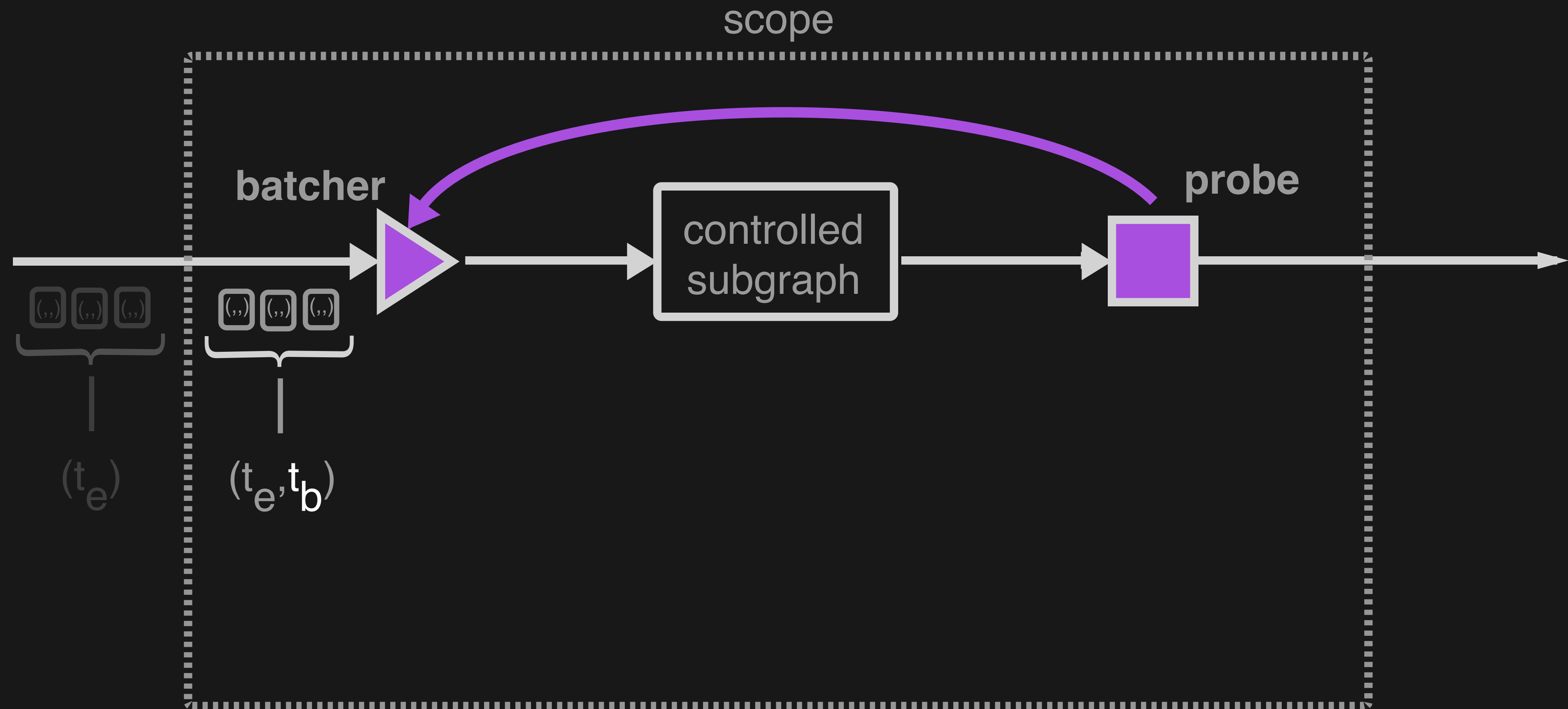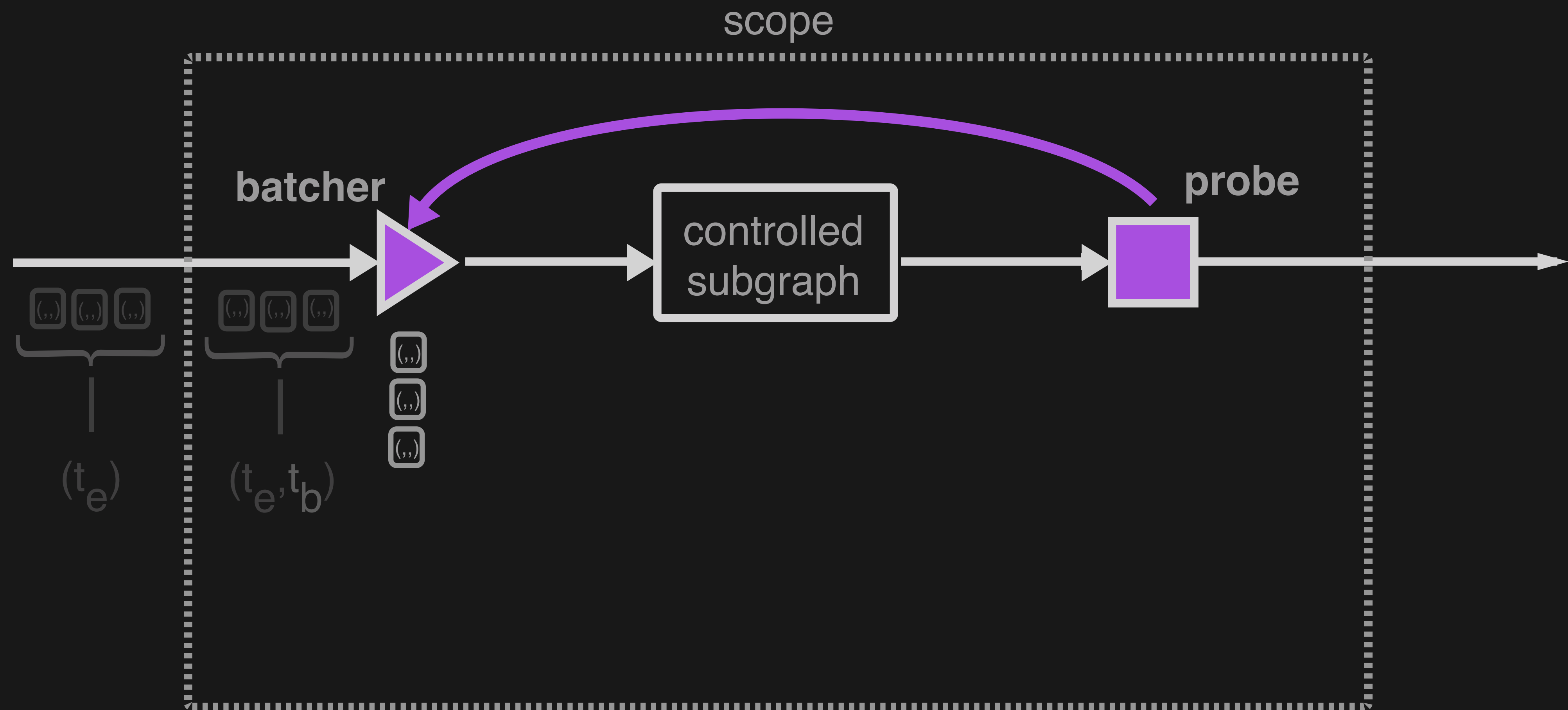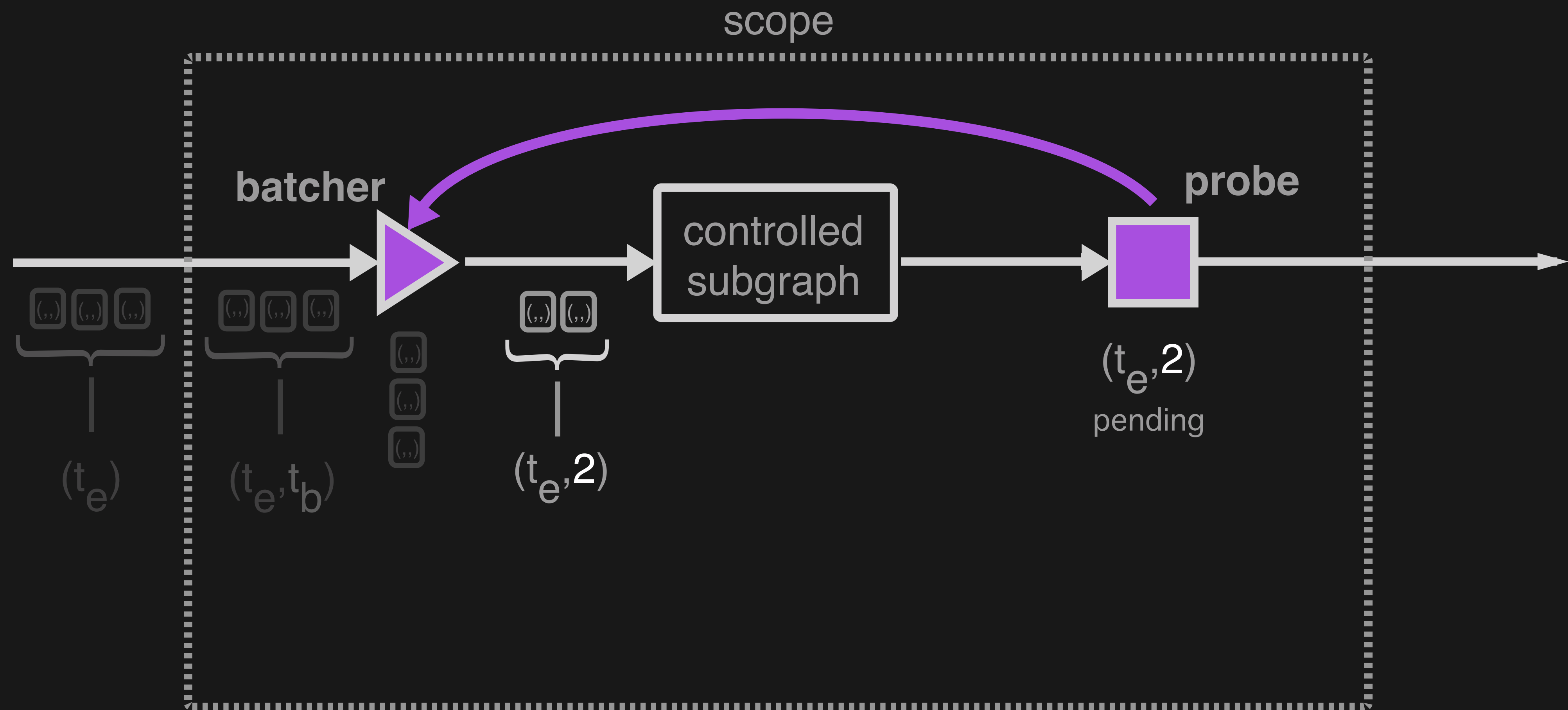
# Faucet - Track batches of *intermediate results*
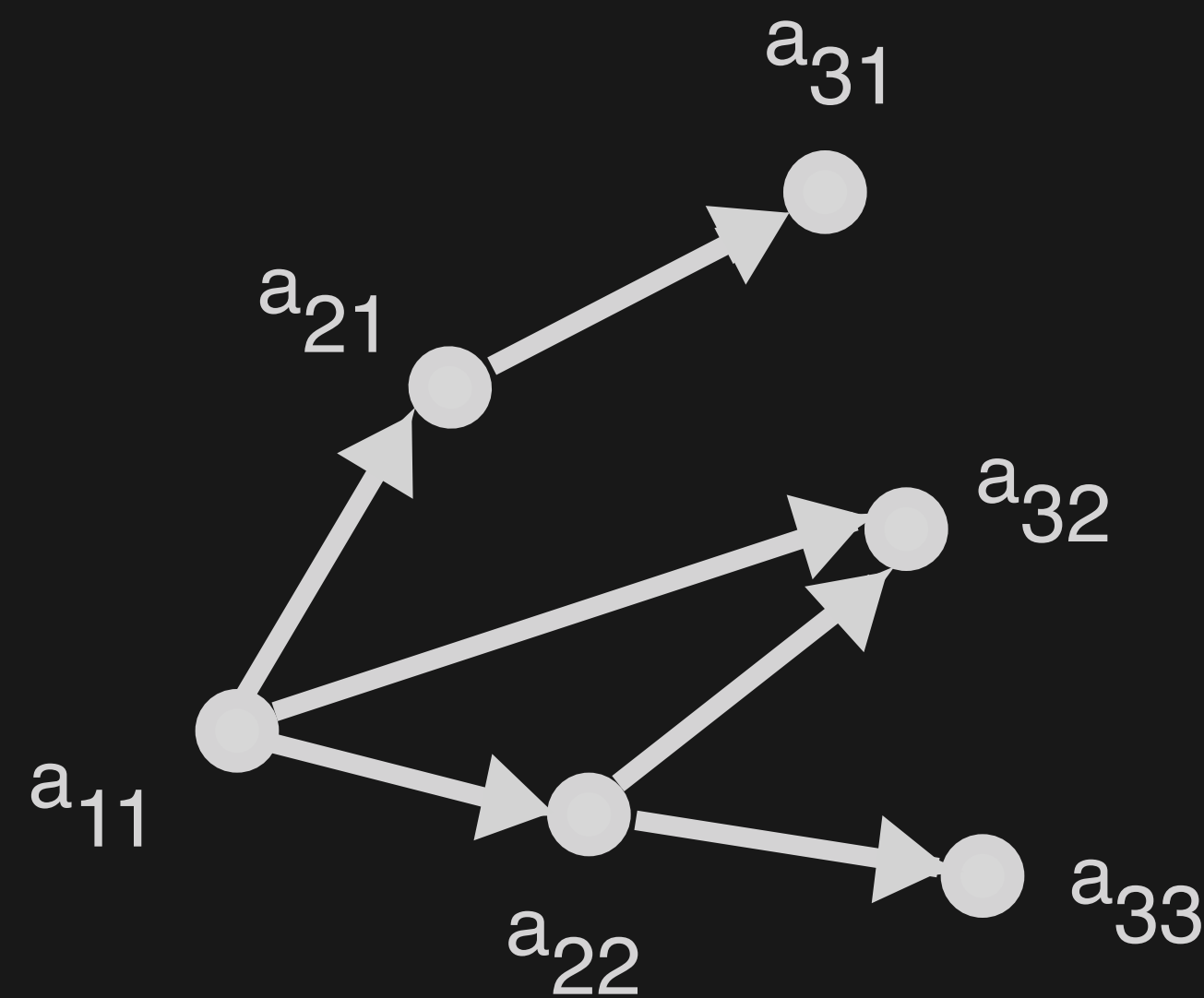
Faucet - Track batches of *intermediate results*

# Example - Enumerate triangles in a directed graph

H. Q. Ngo, C. Ré, and A. Rudra - Generic Join

$a_{31}$

$a_{21}$

$a_{32}$

$a_{11}$
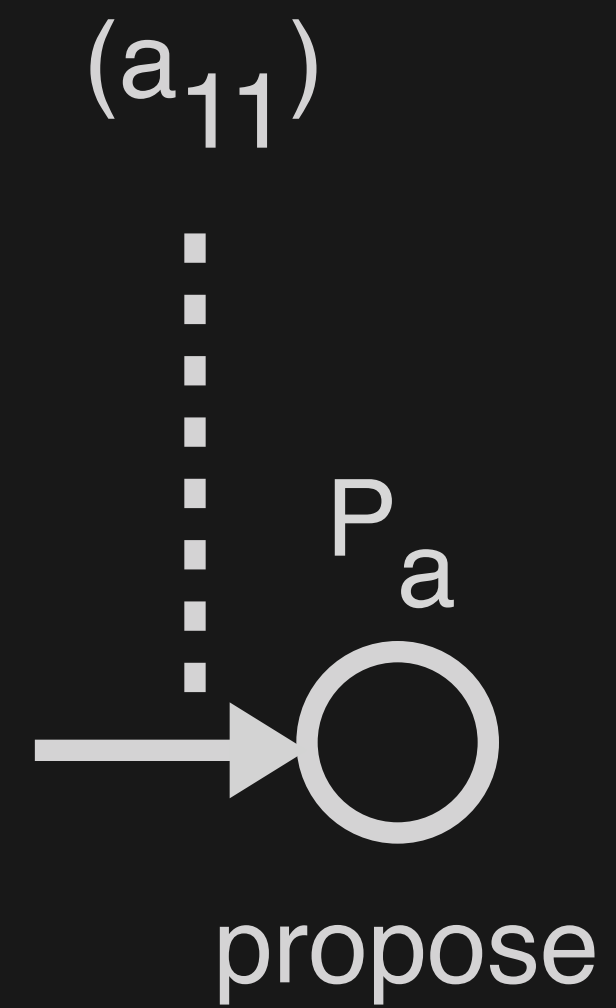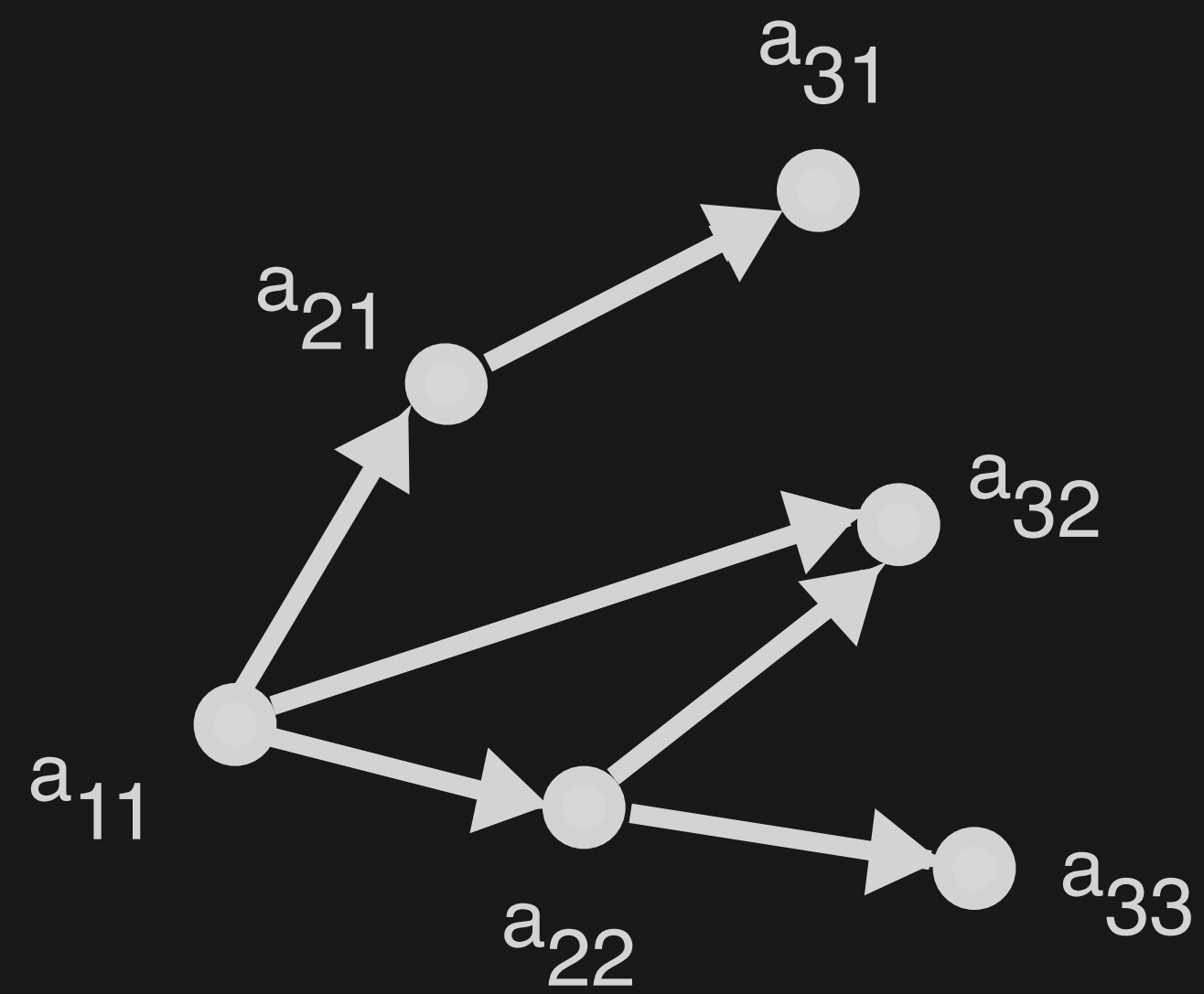
$a_{22}$

$a_{33}$

input graph

build result tuples by extending prefixes

$(a_{11})$

$(a_{11}, a_{22})$

$(a_{11}, a_{22}, a_{32})$

$a_{31}$

$a_{21}$

$a_{32}$

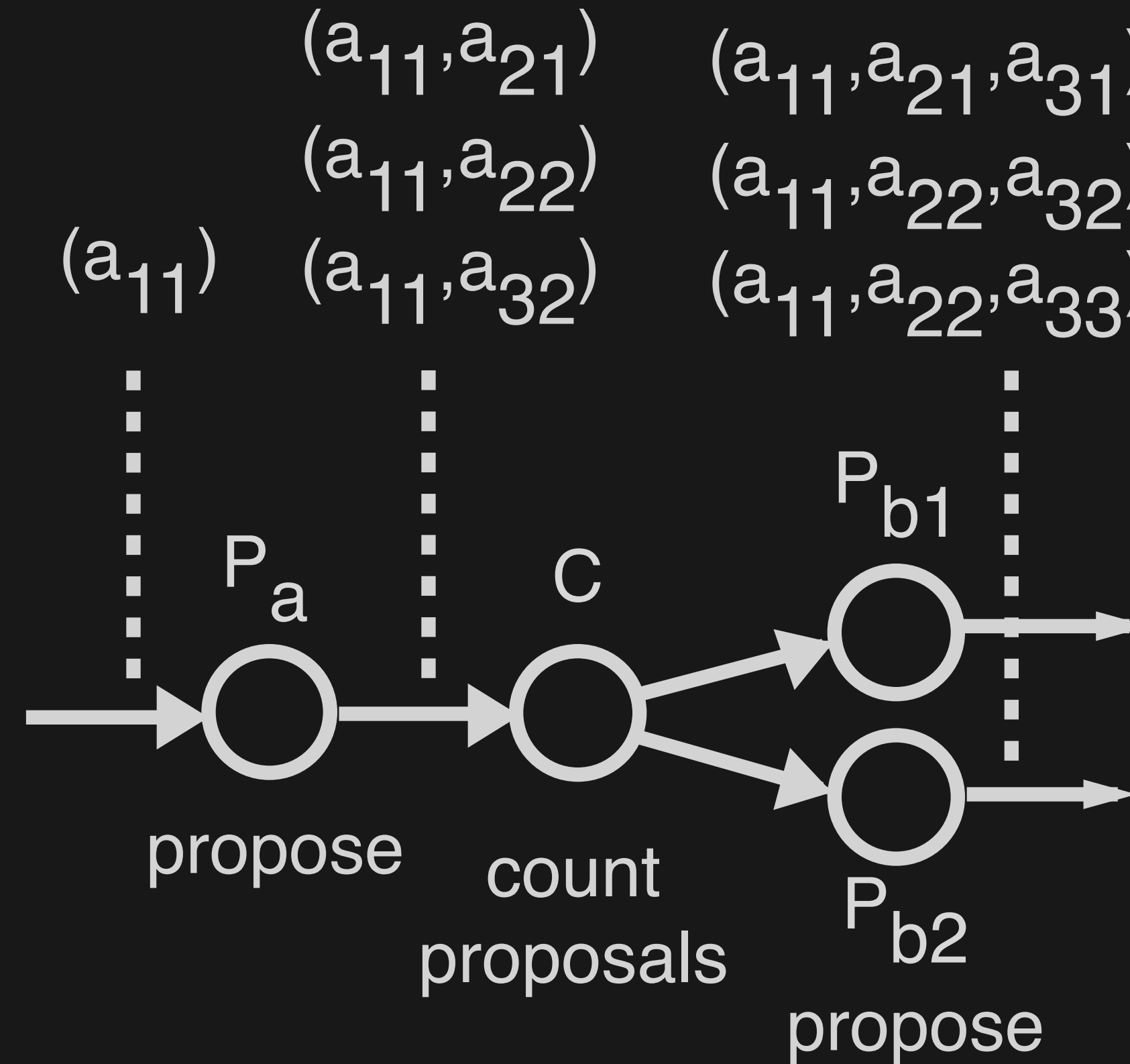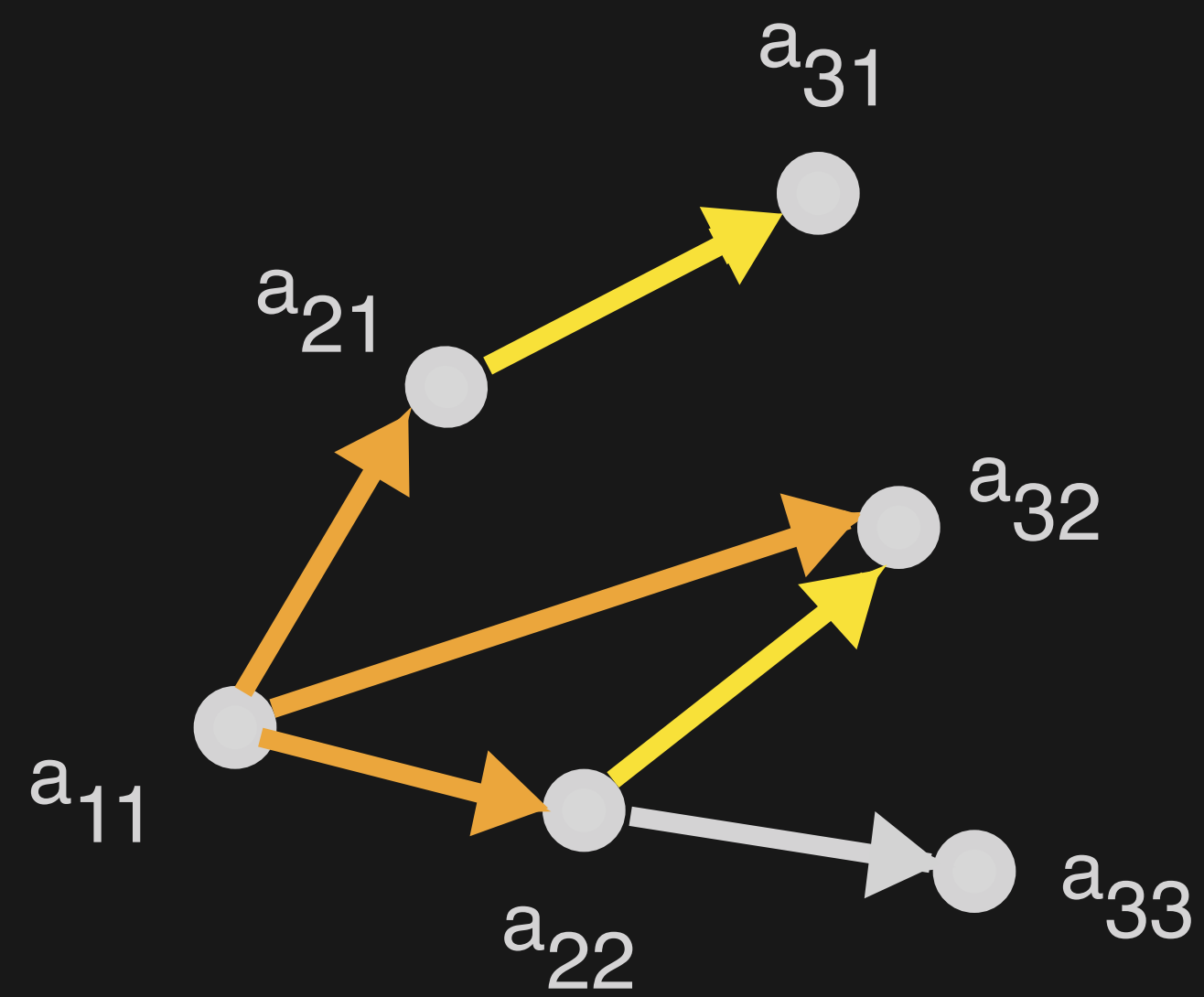$a_{11}$

$a_{22}$

$a_{33}$

$(a_{11})$

$P_a$

propose

# A naïve schedule can generate large *intermediate state*

# A naïve schedule can generate large *intermediate state*

# A naïve schedule can generate large *intermediate state*

# A naïve schedule can generate large *intermediate state*



$(a_{11}, a_{21}, a_{31})$
$(a_{11}, a_{22}, a_{32})$
$(a_{11}, a_{22}, a_{33})$
$(a_{12}, a_{23}, a_{34})$
$(a_{12}, a_{23}, a_{35})$

$(a_{11}, a_{21})$
$(a_{11}, a_{22})$
$(a_{12}, a_{23})$

$(a_{11})$
$(a_{12})$

$P_a$ — propose

$C$ — count proposals

$P_{b1}$
$P_{b2}$ — propose

$I_1$
$I_2$ — intersect

$T$

# Faucet *limits buffered intermediate results*



$(a_{11})$
$(a_{12})$

$(a_{11}, a_{21})$
$(a_{11}, a_{22})$
$(a_{12}, a_{23})$

$(a_{11}, a_{21}, a_{31})$
$(a_{11}, a_{22}, a_{32})$
$(a_{11}, a_{22}, a_{33})$
$(a_{12}, a_{23}, a_{34})$
$(a_{12}, a_{23}, a_{35})$

large
output
rate

11

# Faucet *limits buffered intermediate results*

# Faucet *limits buffered intermediate results*



$(a_{11})$
$(a_{12})$

$(a_{11},a_{21})$
$(a_{11},a_{22})$
$(a_{12},a_{23})$

$(a_{11},a_{21},a_{31})$
$(a_{11},a_{22},a_{32})$
$(a_{11},a_{22},a_{33})$
$(a_{12},a_{23},a_{34})$
$(a_{12},a_{23},a_{35})$

# Faucet *limits buffered intermediate results*

# Evaluation - Dataset

Enumerate triangles in the **Livejournal Dataset**

4'847'571 nodes

68'993'773 edges

285'730'264 triangles

Hardware

Intel Xeon E5-2650 @ 2.00GHz

16 physical cores

10Gbps link

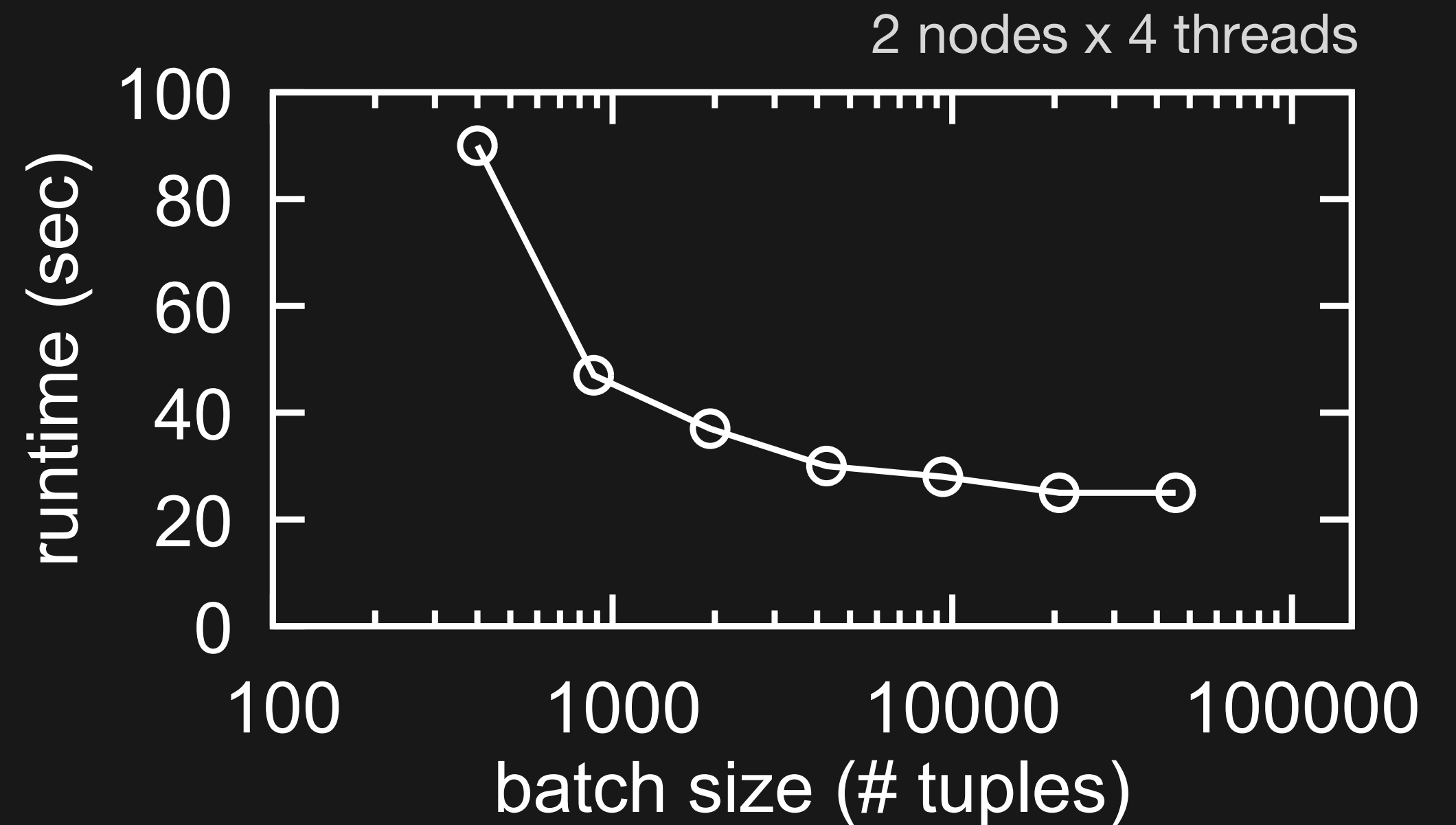# Evaluation - Sensitivity to parameter choice

**N**<sub>**batches**</sub>  number of batches

in-flight in parallel

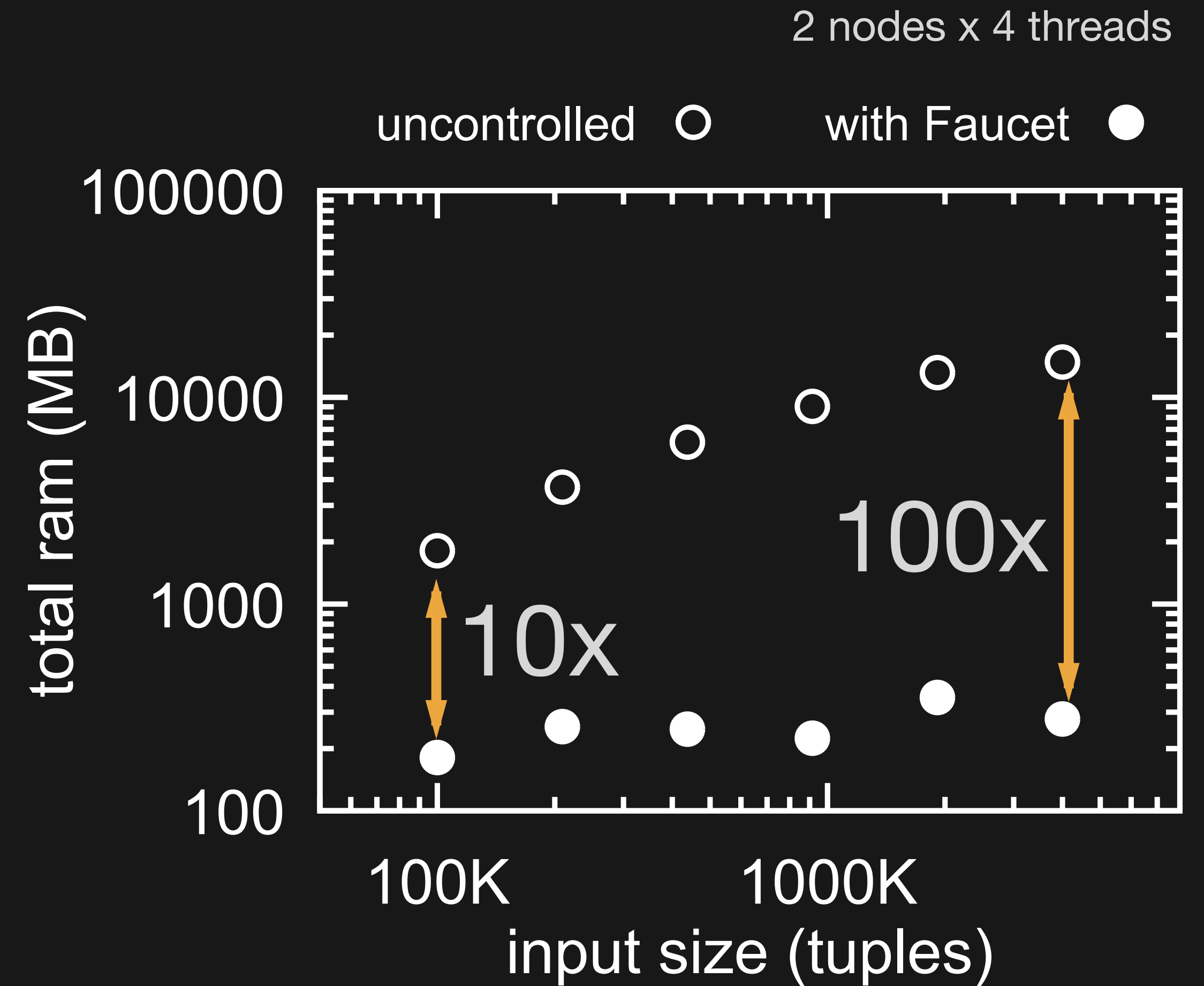$N_{batches} \geq 2$ mitigates stragglers

**B**  batch size

2 nodes x 4 threads



batch size (# tuples)

# Evaluation

## Memory savings



Runtime overhead

**15-25%**

2 nodes x 4 threads

uncontrolled ○     with Faucet ●

100x

10x

total ram (MB): 100000, 10000, 1000, 100

input size (tuples): 100K, 1000K

Faucet

scope

controlled
subgraph

batcher          probe
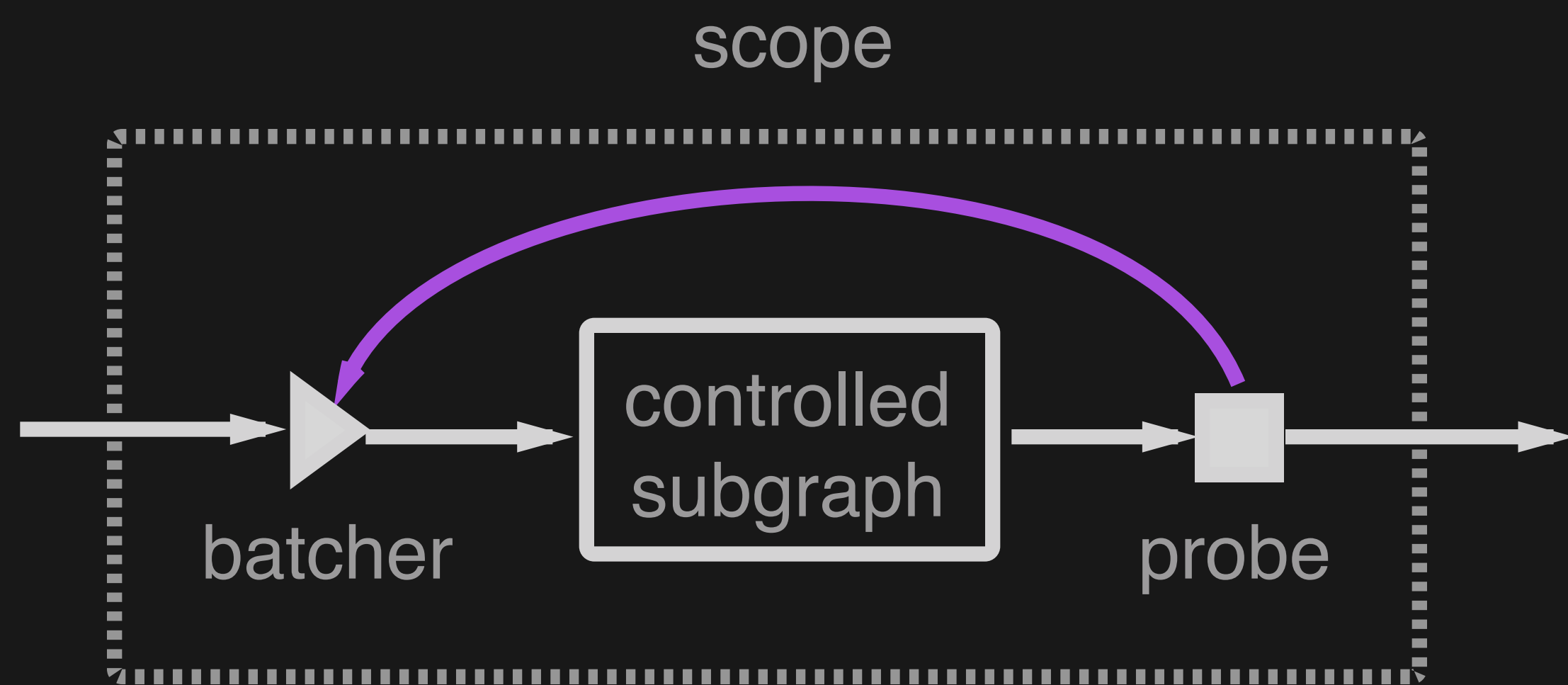
limits intermediate state

RAM is increasingly the main
cost of a system

Memory savings          10-100x
                        or more

Overhead                15-25%

15