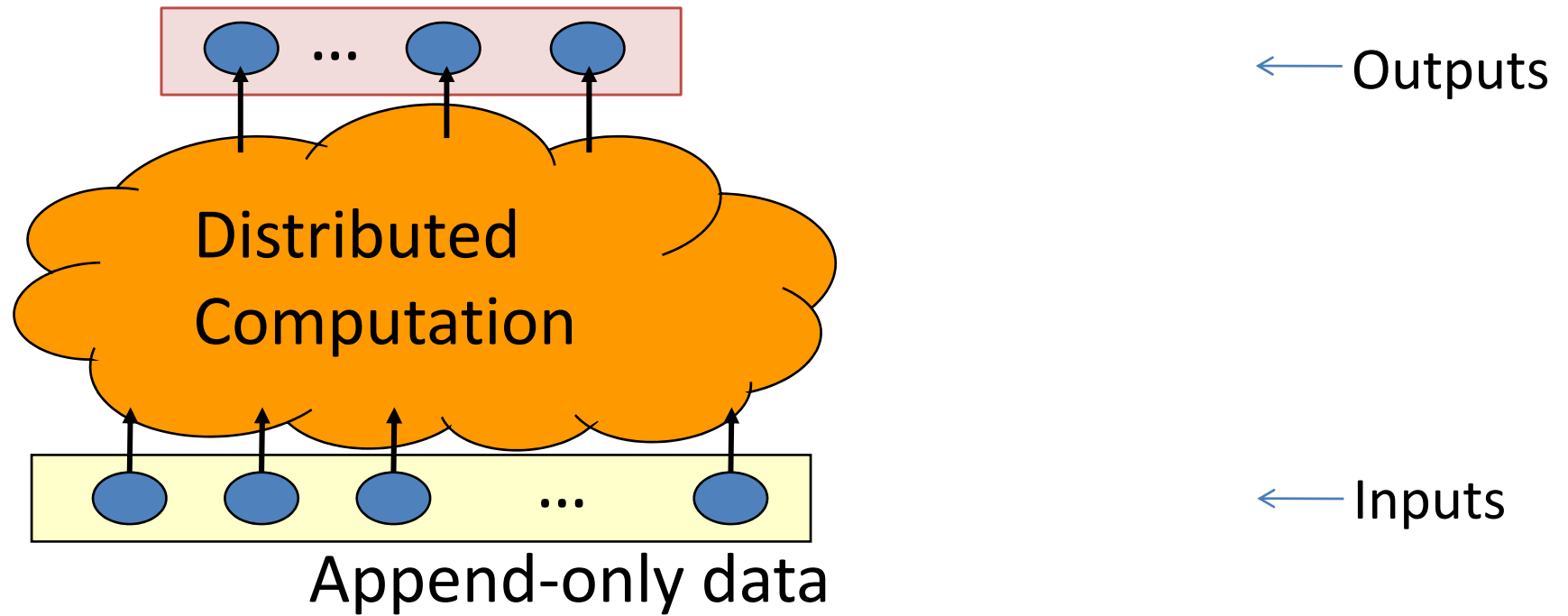# DryadInc: Reusing work in large-scale computations

*Lucian Popa[*+], Mihai Budiu[+],*
*Yuan Yu[+], Michael Isard[+]*
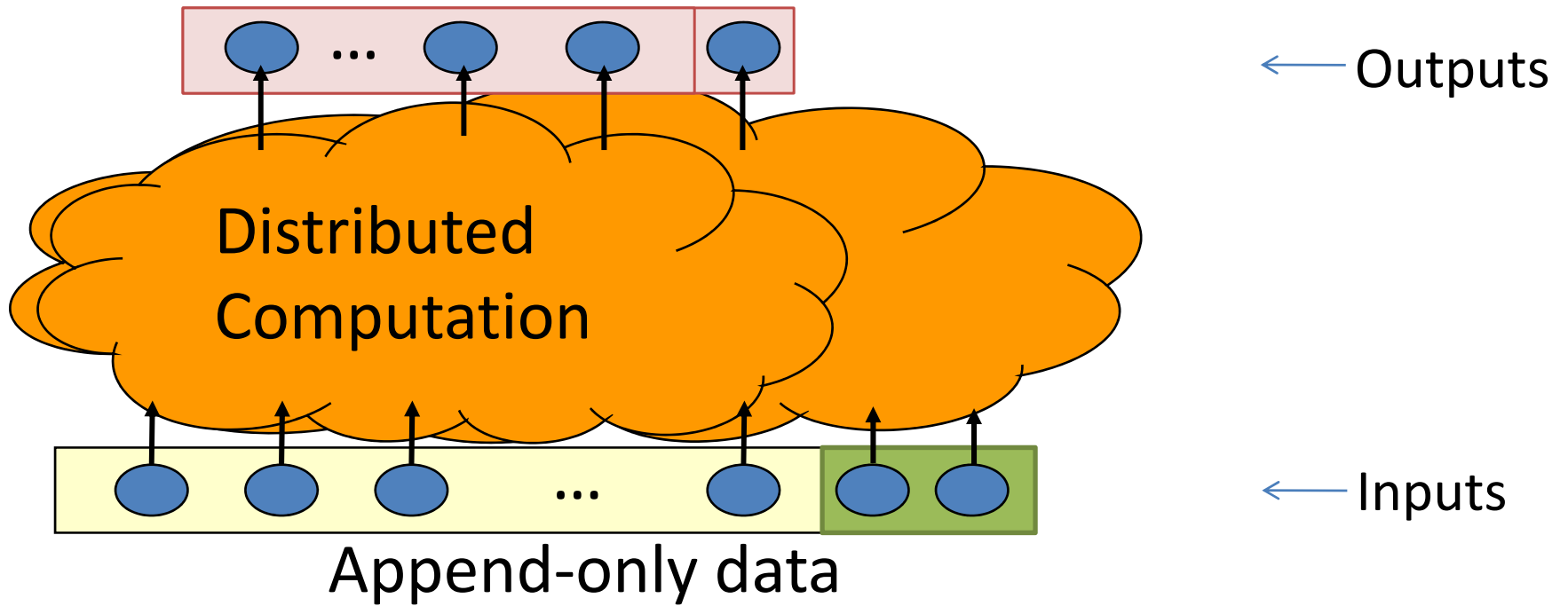
[+] Microsoft Research Silicon Valley
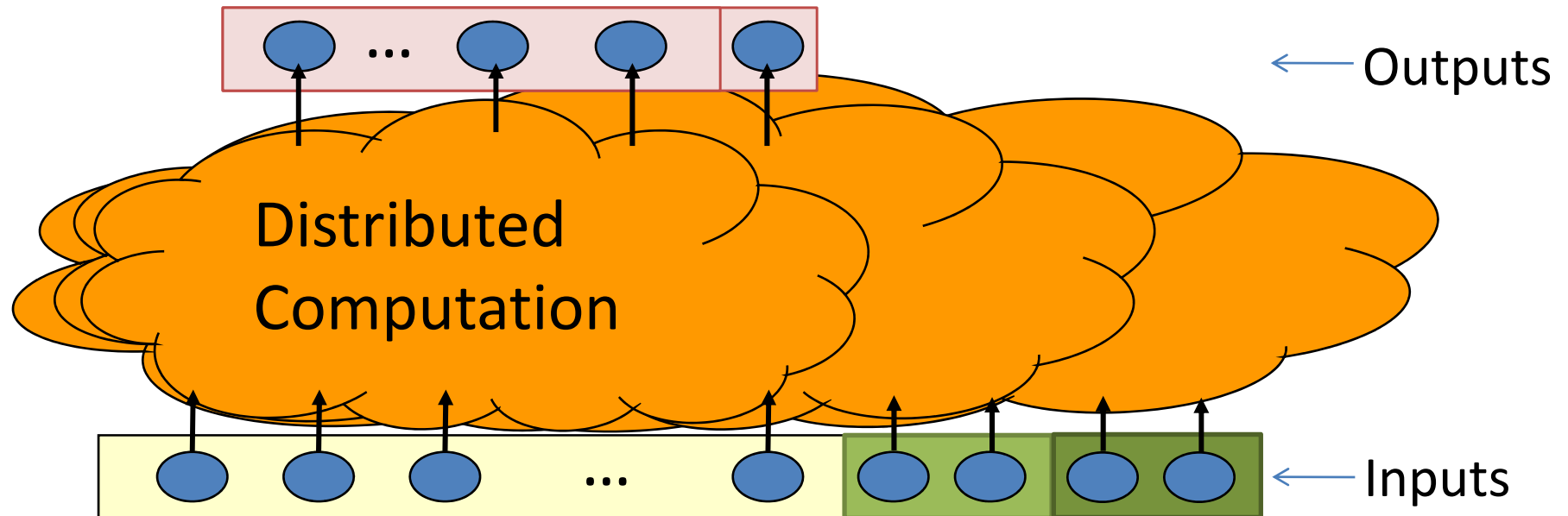[*] UC Berkeley

# Problem Statement



Distributed Computation

Append-only data

← Outputs

← Inputs

# Problem Statement

# Problem Statement



Outputs

Distributed Computation

Inputs

**Goal: *Reuse* (part of) prior computations to:**
- Speed up the current job
- Increase cluster throughput
- Reduce energy and costs

# Propose Two Approaches

**1. IDE**

Reuse *IDEntical computations* from the past
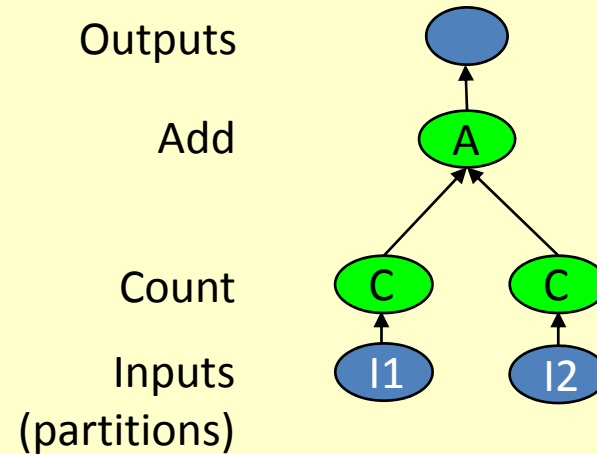
(like `make` or memoization)

**2. MER**

Do only *incremental computation* on the new data

and *MERge* results with the previous ones

(like `patch`)

# Context

- Implemented for **Dryad**
  - Dryad Job = Computational DAG
    - **Vertex:** arbitrary computation + inputs/outputs
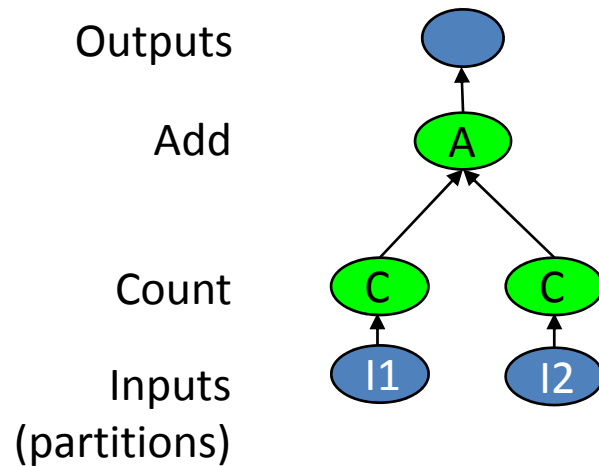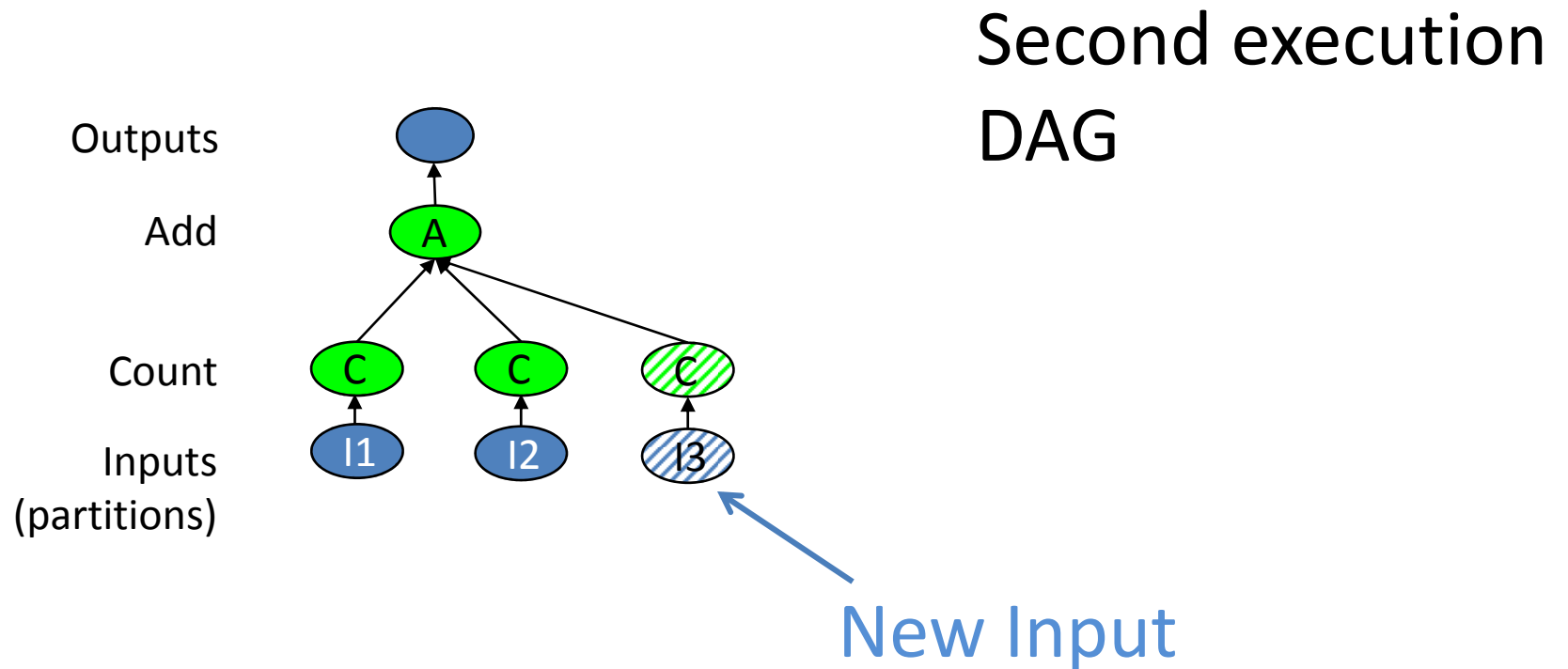    - **Edge:** data flows

Simple Example:
Record Count

Outputs

Add

Count

Inputs
(partitions)

# IDE – IDEntical Computation

Record Count

First execution
DAG



Outputs

Add

Count

Inputs
(partitions)

# IDE – IDEntical Computation

Record Count

Second execution
DAG

Outputs

Add

Count

Inputs
(partitions)

New Input

# IDE – IDEntical Computation

Record Count

Second execution DAG



Outputs

Add

Count

Inputs (partitions)
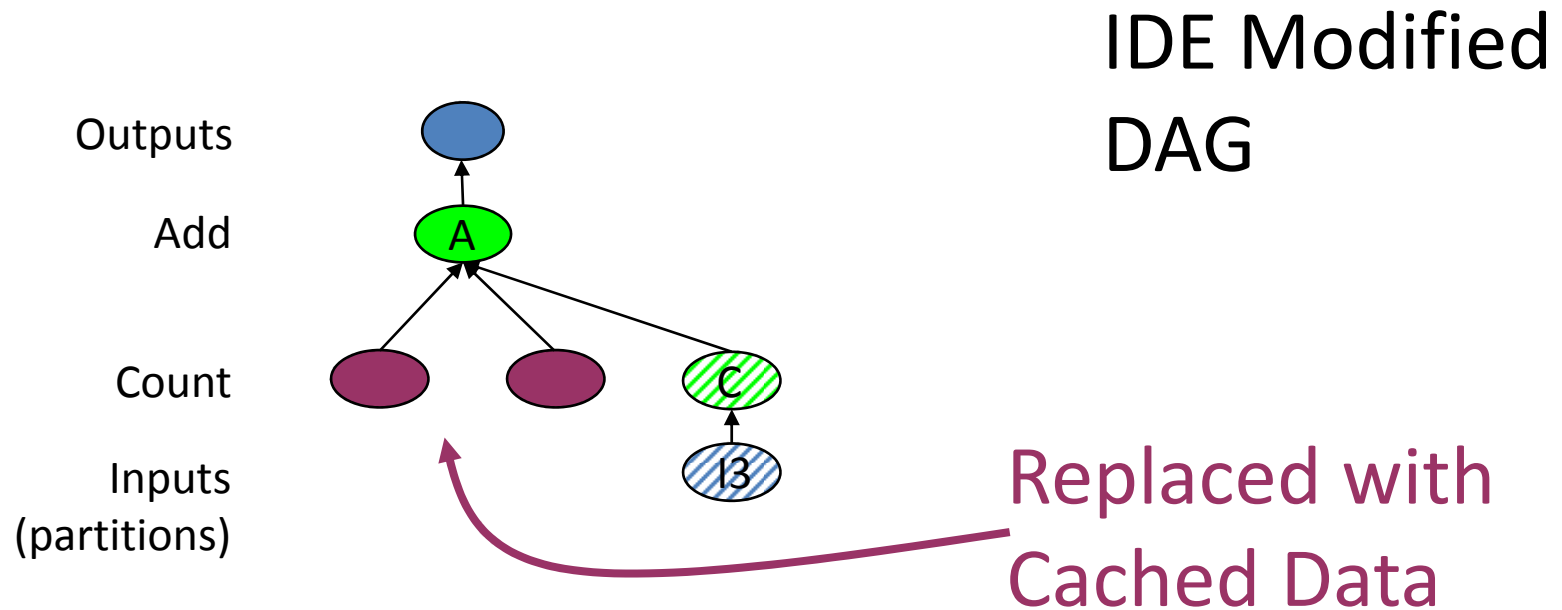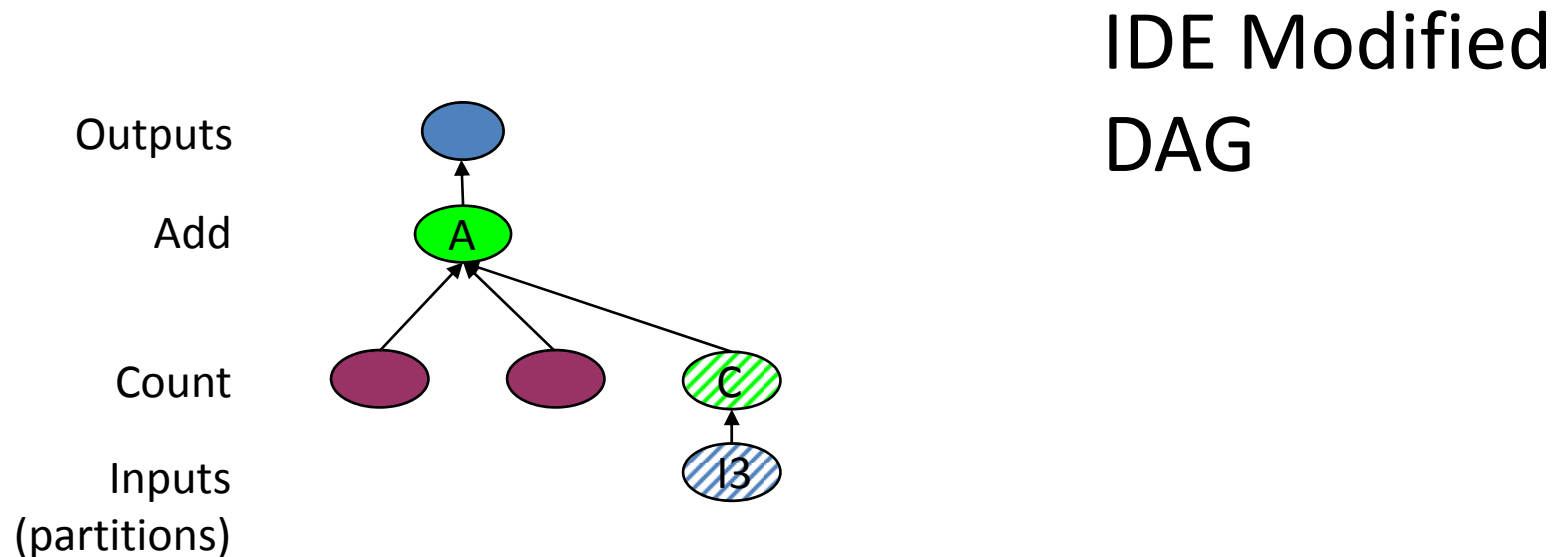
Identical subDAG

# IDE – IDEntical Computation

Replace identical computational subDAG with edge data cached from previous execution



IDE Modified DAG

Outputs

Add    A

Count

Inputs
(partitions)    C

I3

Replaced with Cached Data

# IDE – IDEntical Computation

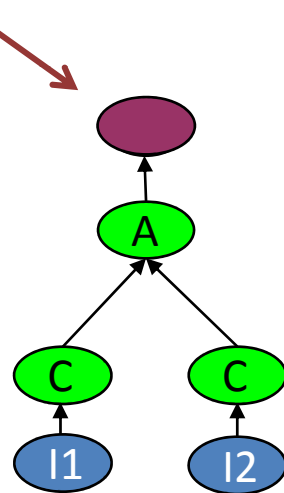Replace identical computational subDAG with edge data cached from previous execution

IDE Modified DAG

Outputs
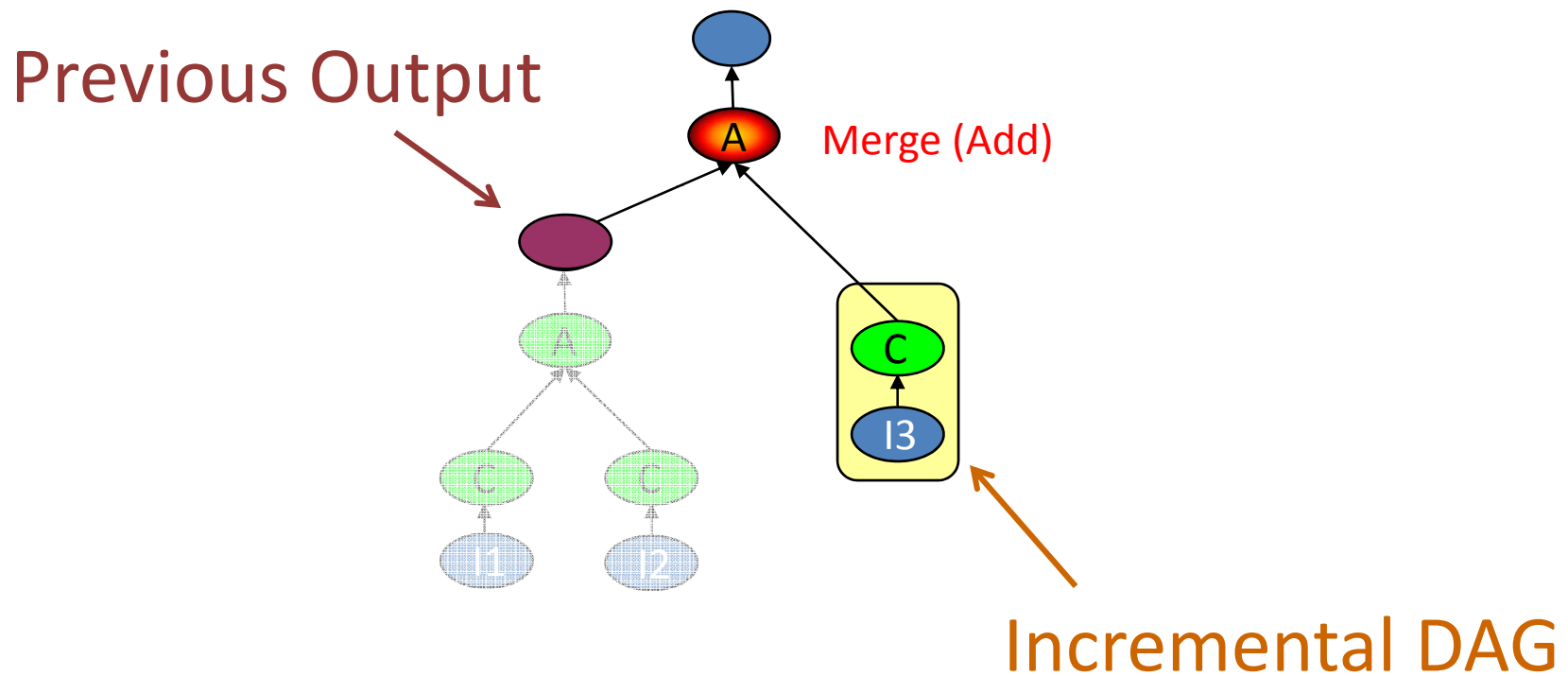
Add

Count

Inputs (partitions)



Use DAG *fingerprints* to determine if computations are identical
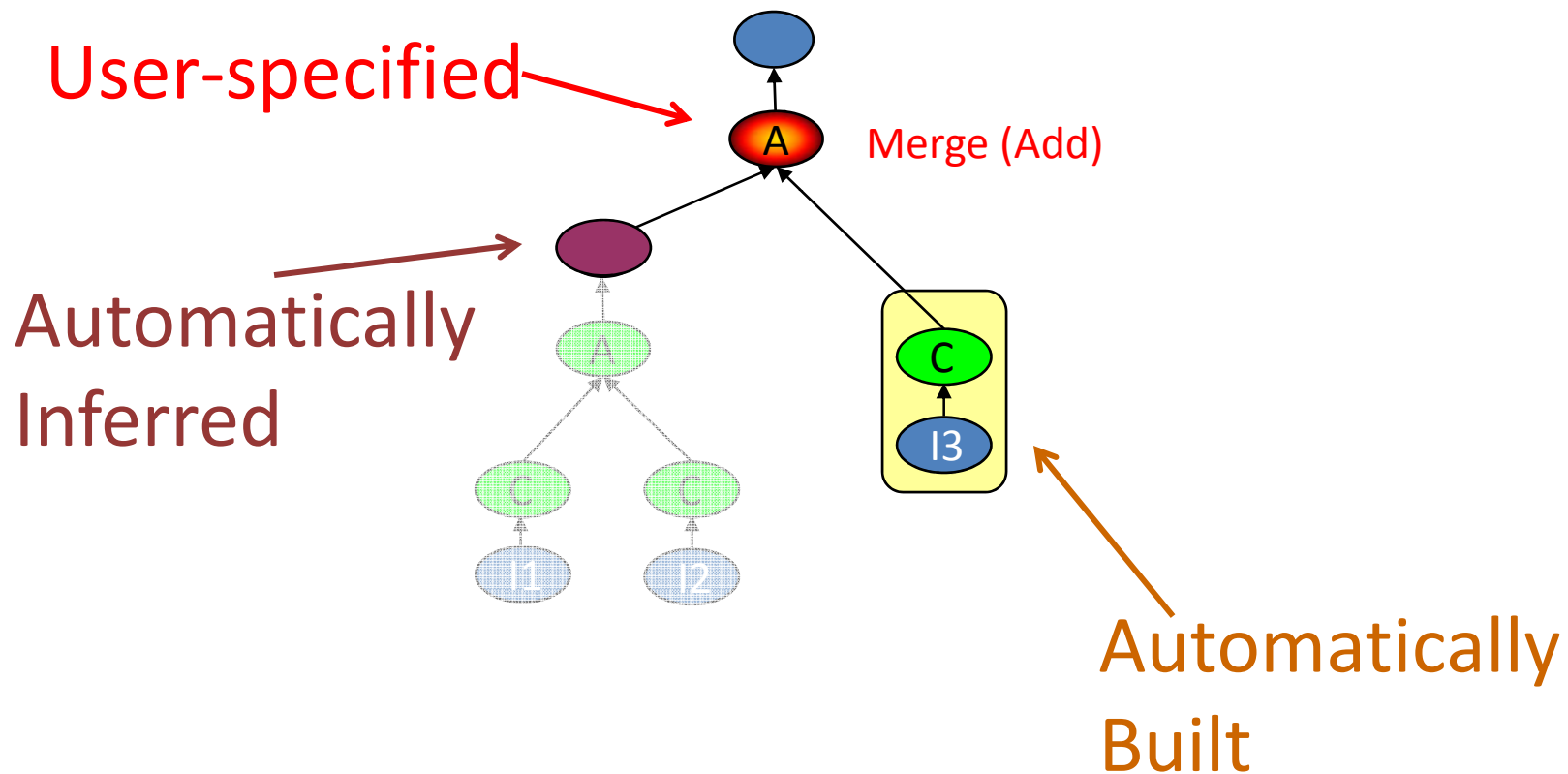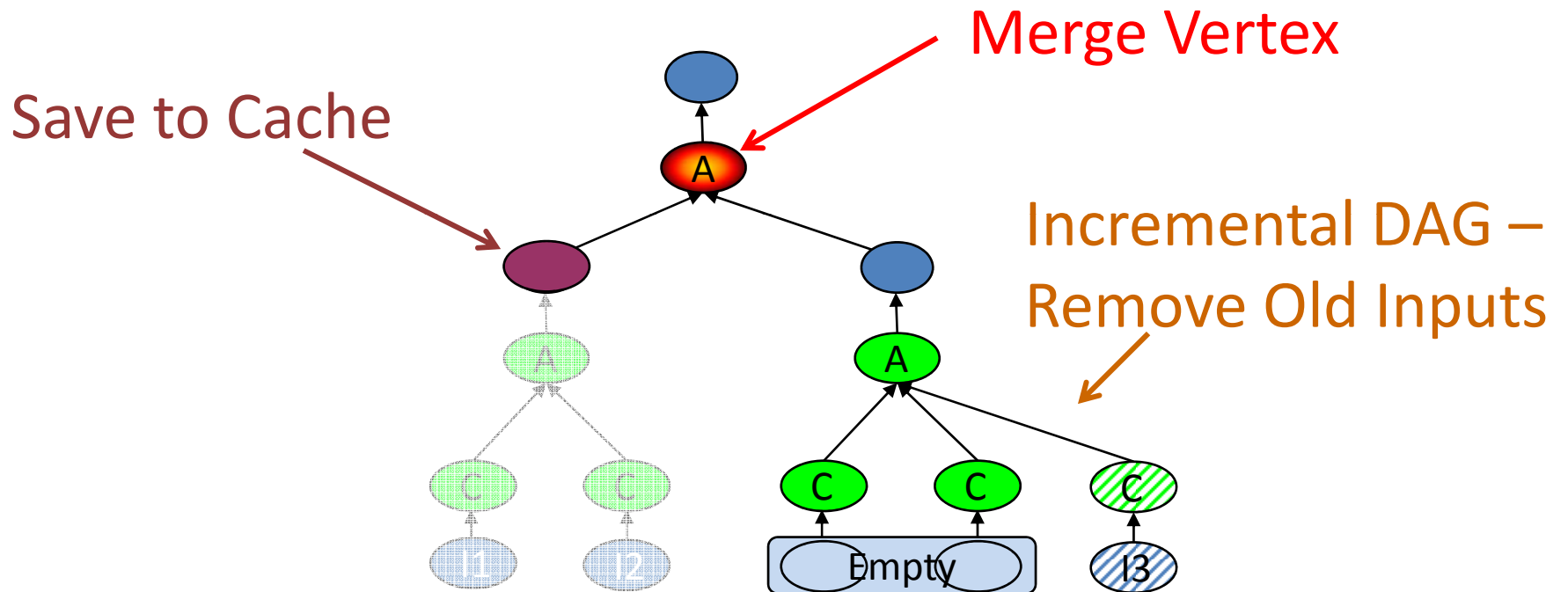
# Semantic Knowledge Can Help

Reuse Output

# Semantic Knowledge Can Help
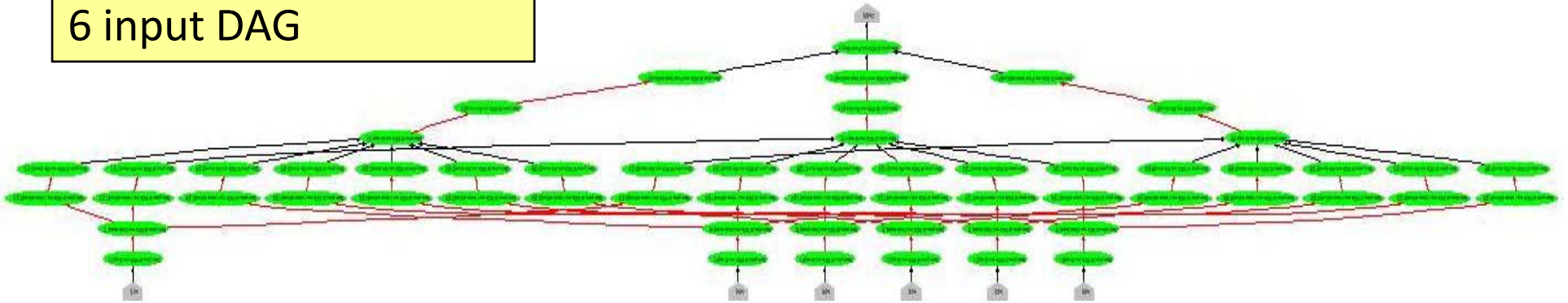
# MER – MERgeable Computation
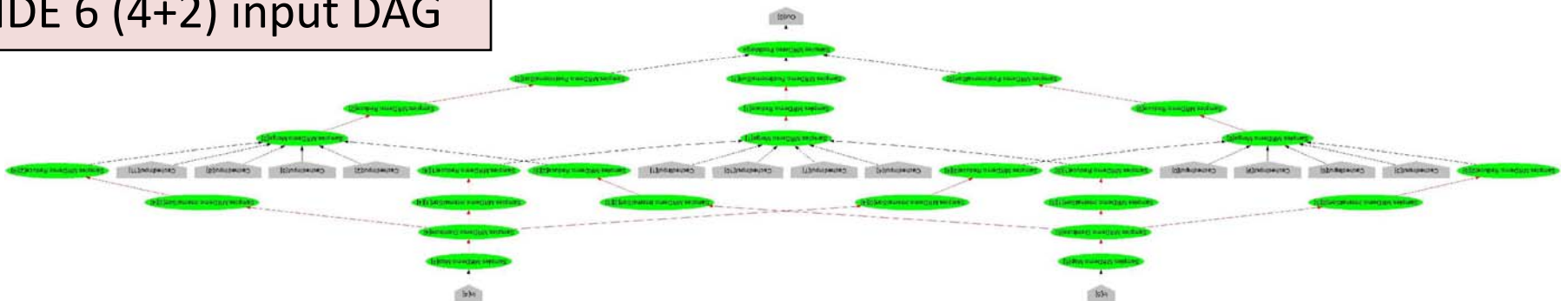
# MER – MERgeable Computation
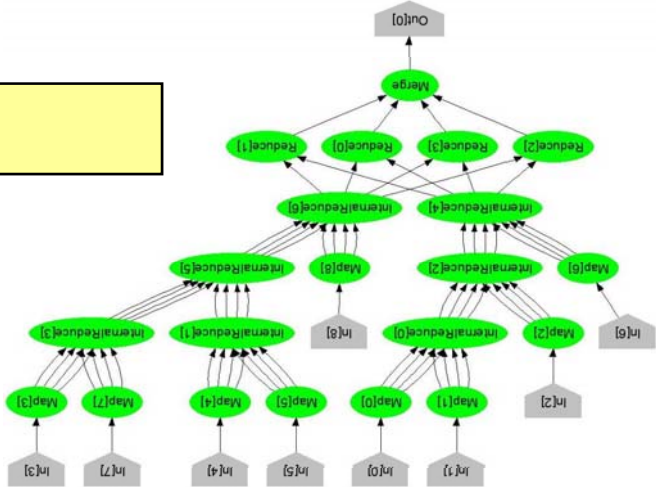
# IDE in practice

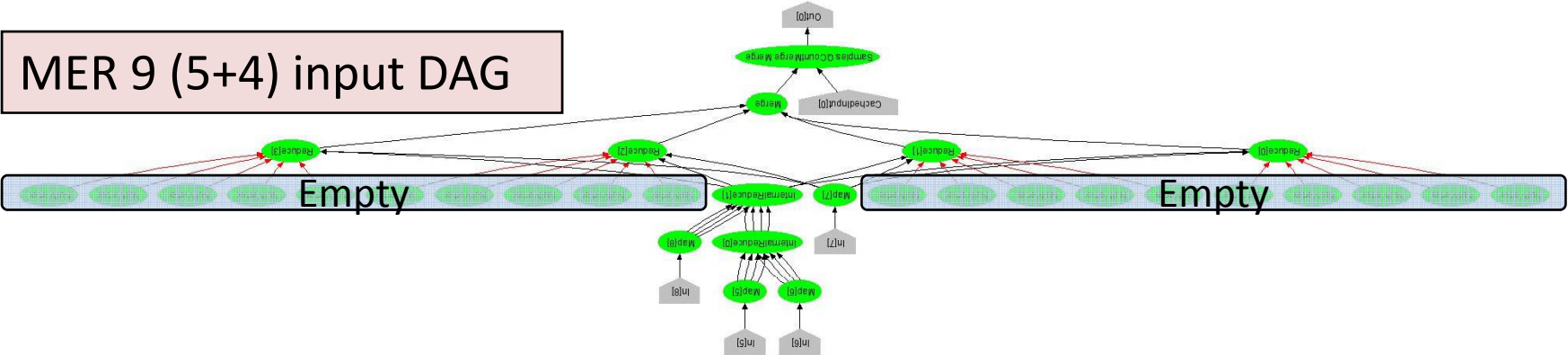6 input DAG

IDE 6 (4+2) input DAG
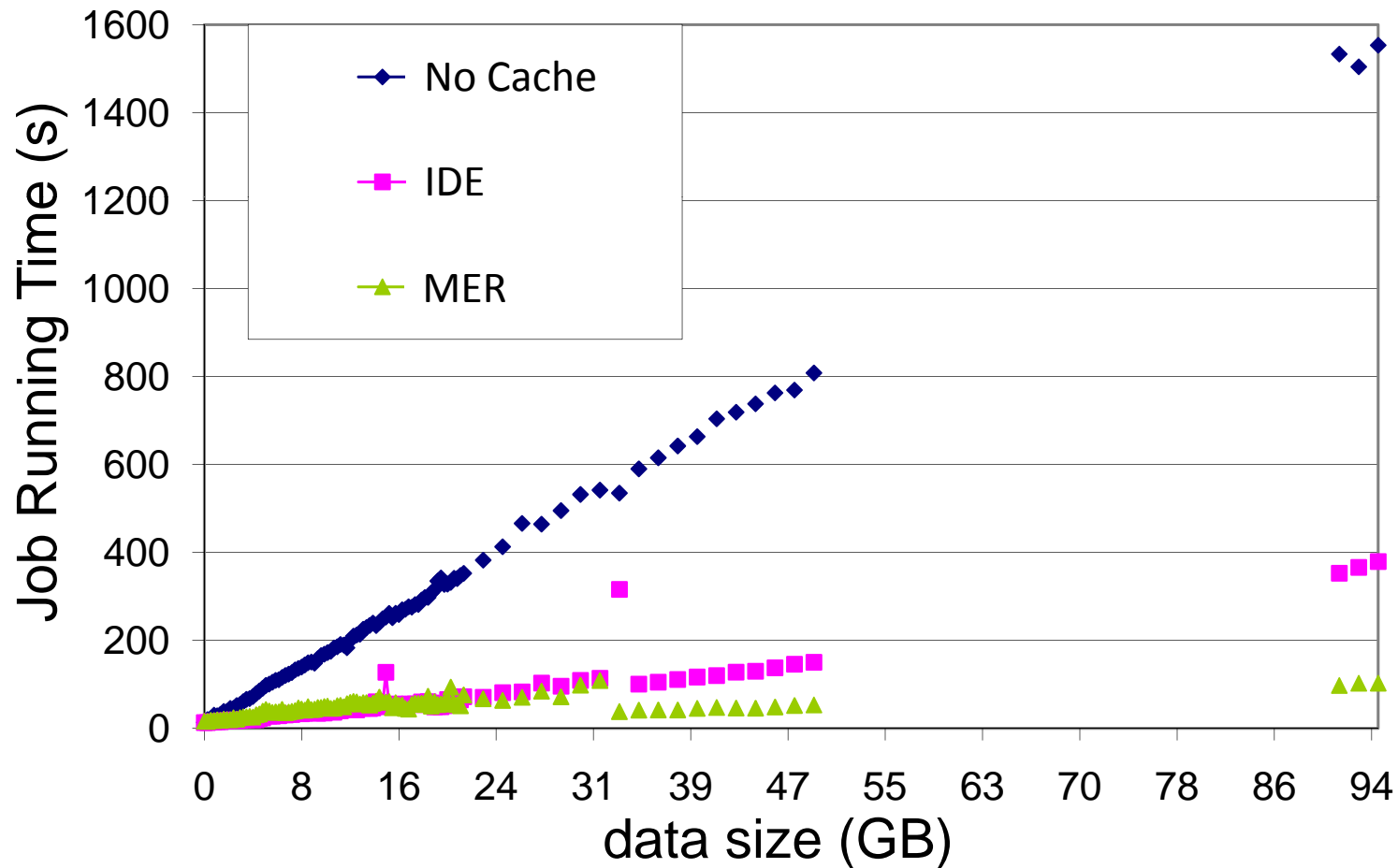
# MER in practice

9 input DAG

MER 9 (5+4) input DAG

# Evaluation – Running time

## Word Histogram Application – 8 nodes

# Discussion

- **MapReduce:** just a particular case
  - IDE reuses the output of Mappers
  - MER requires combined Reduce function

- **Combine IDE with MER:** benefits don't add up
  - IDE can be used for the incremental DAG at MER

- **More semantic knowledge:** further opportunities
  - Generate merge function automatically
  - Improve incremental DAG

- **Sliding window on input data:** IDE works unchanged, MER requires "divide" besides merge

# Conclusions & Questions

- **Problem:** reuse work in distributed computations on append-only data

- **Two methods:**
  - **INC** – reuse IDEntical past computations
    - No user effort
  - **MER** – MERge past results with new ones
    - Small user effort, potentially larger gains

- Implemented for **Dryad**