# Blaze:
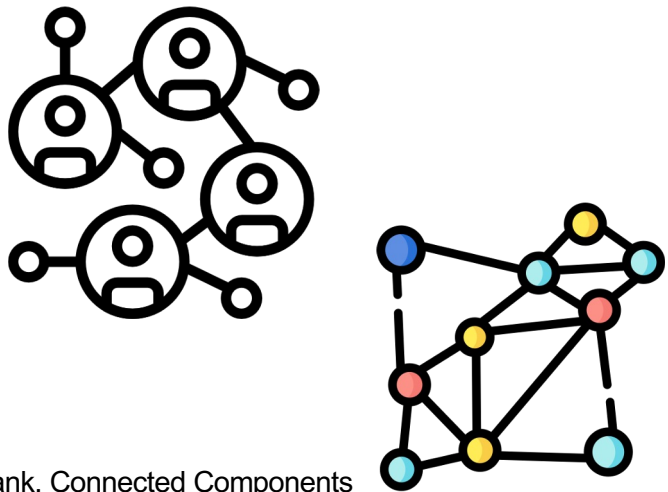# Holistic Caching for Iterative Data Processing

**Won Wook SONG**, Jeongyoon Eo, Taegeon Um, Myeongjae Jeon, Byung-Gon Chun

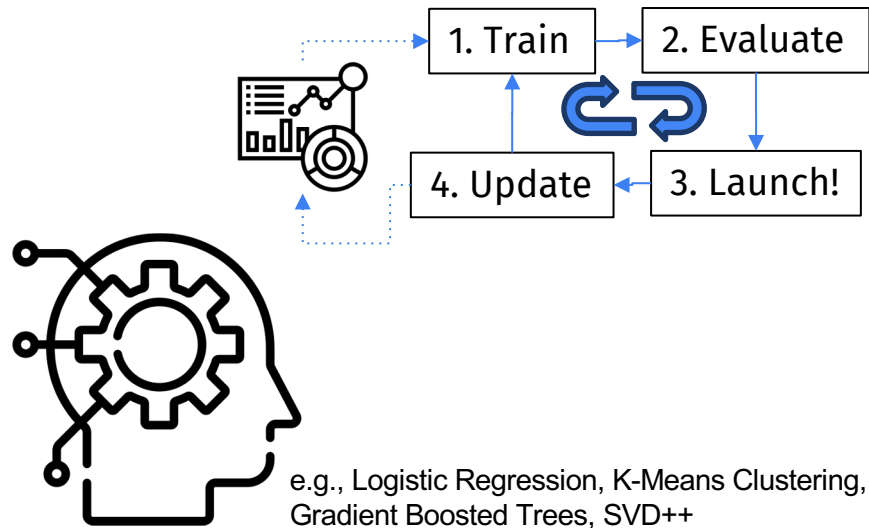Seoul National University, Samsung Research, UNIST, FriendliAI

# Caching is Essential in Complex Iterative Data Processing

## Graph Processing



e.g., PageRank, Connected Components

## Iterative Machine Learning



1. Train → 2. Evaluate → 3. Launch! → 4. Update → (loop)

e.g., Logistic Regression, K-Means Clustering, Gradient Boosted Trees, SVD++

Without caching, intermediate data must be recomputed every time by default

# Cache Hints for Indicating Data to Cache

```
1   // Initialize the PageRank graph
2   val rankGraph = graph.join(...).map(...)
3   for (until convergence) {
4     rankGraph.cache()
5     val rankUpdates = rankGraph.aggregate(...)
6     prevRankGraph = rankGraph
7     rankGraph = rankGraph
8           .outerJoinVertices(rankUpdates) {...}
9     // Submit job and materialize vertices
10    rankGraph.edges.foreachPartition(...)
11    prevRankGraph.unpersist()
12  }
```

Caching decisions are made by users on an operator dataset level

→ **Caches all partitions** of the dataset

# Caching is Done Repeatedly over Iterations until Convergence
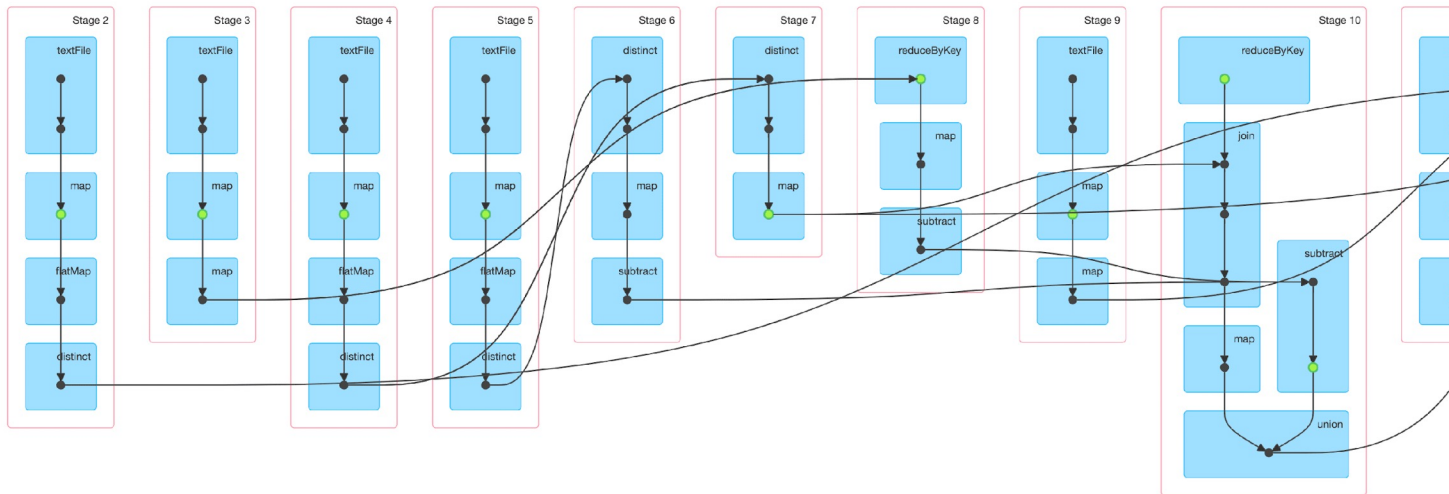
```
1   // Initialize the PageRank graph
2   val rankGraph = graph.join(...).map(...)
3   for (until convergence) {
4     rankGraph.cache()
5     val rankUpdates = rankGraph.aggregate(...)
6     prevRankGraph = rankGraph
7     rankGraph = rankGraph
8             .outerJoinVertices(rankUpdates) {...}
9     // Submit job and materialize vertices
10    rankGraph.edges.foreachPartition(...)
11    prevRankGraph.unpersist()
12  }
```

The complete job DAG is unknown
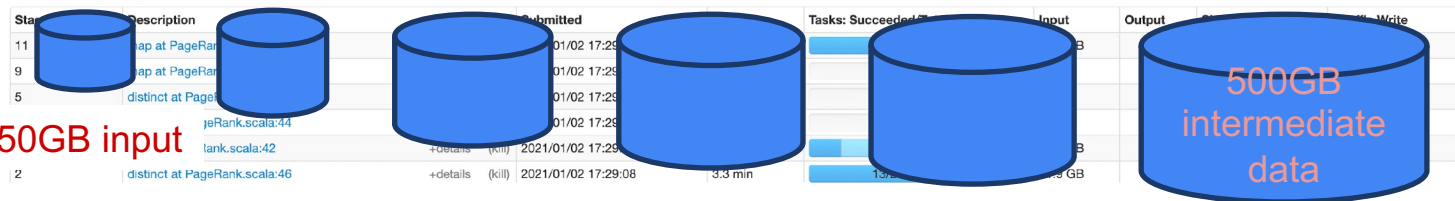until convergence
(one job == one iteration)

→ **Caching is done in a greedy fashion, due to unknown future data references and dependencies**
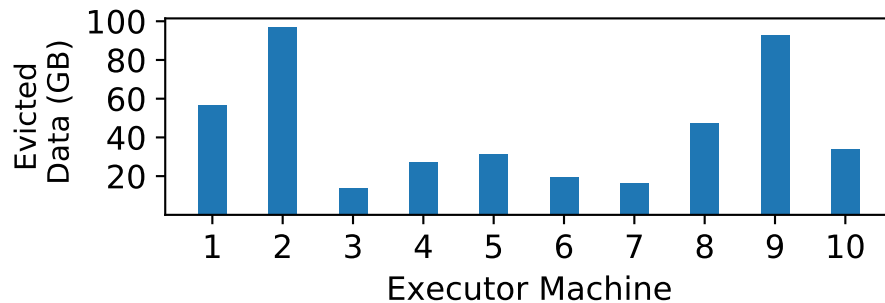
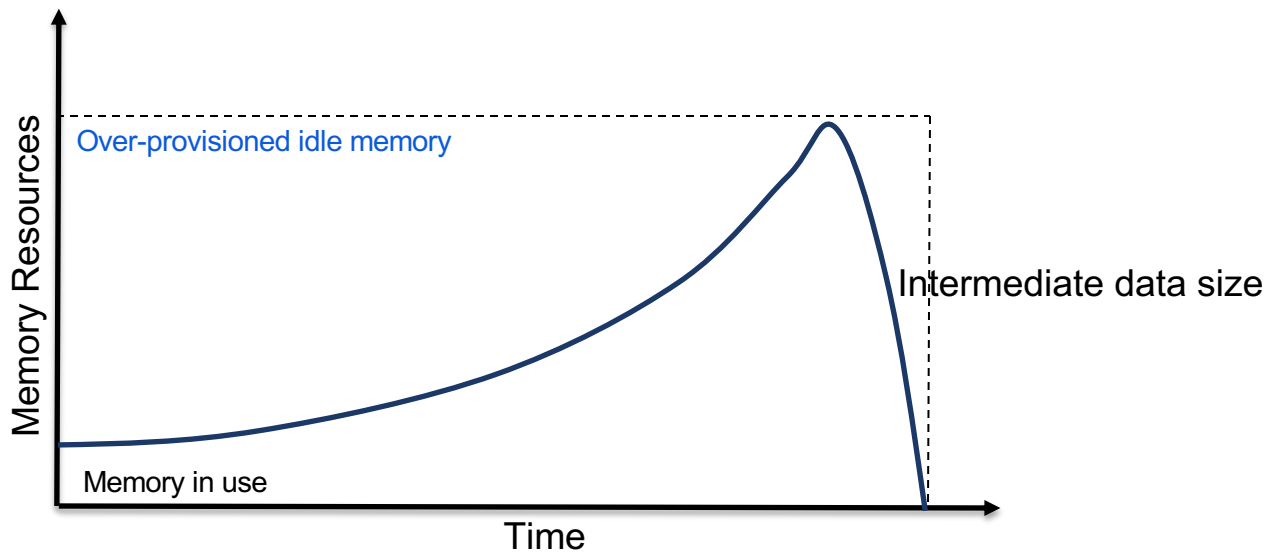# Cache Data Size Increases Over Iterations

# Cache Data Size is Inconsistent Among Executors



Inconsistent memory usage → Bottleneck executors make optimizations tricky
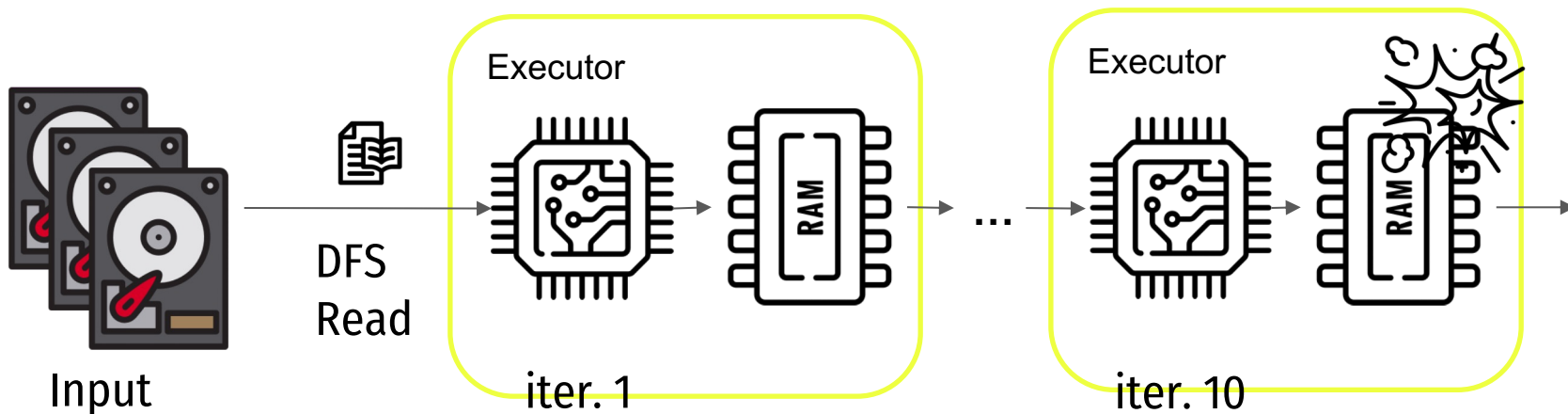
# A Naïve Solution: Over-provisioning Memory



Simplest solution, but reserving memory of *>10x the input size* is *costly*

Plus, we cannot *foresee* when the iterations will *conclude*

# Problem: Memory Space is Constrained



Input → DFS Read → Executor (iter. 1) → ... → Executor (iter. 10)

Also, caching in memory is difficult to scale

(Need to increase # of VMs)

# Goal: Reducing Recomputation Overhead



Recomputation overhead is nontrivial

# Why is this a problem in real-world environments?



Recomputation costs increase exponentially across iterations upon cache recovery

# Using Disks as Secondary Caching Storages



Data can be evicted from memory and spilled on disk, if the memory is full

# Goal: Reducing Disk I/O Overhead



But incurs serialization/deserialization overheads and disk read/write overheads,

which can be larger than recomputation overheads

# Recomputation vs. Disk I/O Overhead



Recomputation overheads increase in latter iterations,

Disk I/O overheads increase with larger intermediate data,

→ Recomputation cost and disk I/O costs *differ* according to application

# Separate Caching/Eviction/Recovery Layers



Caching in existing systems occur in three separate operational layers:

1) **caching** layer, 2) **eviction** layer, 3) **recovery** layer

# Our Goal

## Existing Works

- **Separated** caching mechanism

- **Greedy** cache management

- Based on **pre-defined rules**

- Heuristics based on **past usages**

- **User annotation**-based caching in **dataset** level

## Blaze

- **Unified** caching mechanism

- Derive a **blueprint** of the optimal cache state **in advance**

- Based on **dynamically** measured **metrics**

- **Cost predictions** based on **current** and **future** usages of **cached data partitions**

- **Automatic caching** in **partition** level within the data processing system

# To Cache, or Not To Cache?

**Memory**

A  B  C

D

**Disk**

Whether or not to cache D if memory is full?

- Default action: Spark **always caches D** according to **user annotation**
  → May lead to **unnecessary** caching

- Proposed action: if D incurs smaller overhead compared to A, B, C, we **shouldn't cache D** in the first place

# To Evict, or Not To Evict?

Which partition to evict if memory is full?

- Default action: Spark always evicts partitions based on the **eviction policy** (e.g., LRU)

- Proposed action: evict the partition with the **smallest potential cost** = min(recomputation cost, disk cost)

  - Each partition has different size and potential cost

  - Example:
    - Recomp(A): 3s, Disk(A): 5s, Size(A): 500MB
    - Recomp(B): 5s, Disk(B): 2s, Size(B): 200MB
    - **Recomp(C): 1s**, Disk(C): 5s, Size(C): 500MB
      → Evict C!

# Data Dependency Changes Dynamically and Unpredictably



Recomputation cost for partition B = Computation cost [A → B]

# Data Dependency Changes Dynamically and Unpredictably



Recomputation cost for partition B = Computation cost [Input → A → B]

→ Recomputation cost/length **varies** on the cached partitions

# Data Dependency Changes Dynamically and Unpredictably



Cached partition B is referenced once, by E

Cached partition C is referenced twice, by D and F

# Data Dependency Changes Dynamically and Unpredictably



Cached partition B is referenced 3 times, by D, F (through C), and E

→ # of cache references varies on the cached data

# Blaze Design Principles

1. Workload **DAG profiling** and **dynamic** metric measurements to induce

   **potential costs** and references of each partitions

2. Incorporating **potential** data caching efficiency into a **unified cost**

   **optimization function**

3. **Automatic caching decisions** based on fast ILPs

# System Overview



(1) Execution with Partial Data

Extract future dependencies

Blaze Runner (script)

Spark App

Profiler

Spark Master

**CostLineage**

# System Overview



(1) Execution with Partial Data

Extract **future** dependencies

- Tracking **runtime metrics** of partitions (e.g., which data are stored, the size of data)
- **Estimate costs** with data dependencies

Blaze Runner (script)

Spark App

Profiler

Spark Master

**CostLineage**

(2) Execution with Full Input Data

# System Overview

**Extract future dependencies**

(1) Execution with Partial Data

- Tracking runtime metrics of partitions (e.g., which data are stored, the size of data)
- Estimate costs with data dependencies

Blaze Runner (script)

Spark App

(2) Execution with Full Input Data

Profiler

Spark Master

**CostLineage**

Caching/Eviction Decisions

Spark Executor

Task  Task
Execution Memory
Caching Memory
Caching disk

Spark Executor

Task  Task
Execution Memory
Caching Memory
Caching disk

Spark Executor

Task  Task
Execution Memory
Caching Memory
Caching disk

# Collecting Potential Future References

## Submitted jobs (profiling phase)

## Cost Lineage

S1

| RDD1 |
| RDD2 |
| RDD5 |

S2

| RDD1 |
| RDD2 |
| RDD11 |

S4

| RDD1 |
| RDD2 |
| RDD17 |
| RDD20 |

| RDD21 |
| RDD27 |
| RDD31 |
| RDD40 |

S3

| RDD6 |
| RDD19 |

# Collecting Potential Future References

## Submitted jobs (profiling phase)

S1

S2

S4

```
S1          S2              S4
┌─────────┐ ┌─────────┐  ┌──────────────────────┐
│ RDD1    │ │ RDD1    │  │ RDD1                 │
│  ↓      │ │  ↓      │  │  ↓                   │
│ RDD2    │ │ RDD2    │  │ RDD2                 │
│  ↓      │ │  ↓      │  │  ↓                   │
│ RDD5    │ │ RDD11   │  │ RDD17 ──→ RDD21      │
└─────────┘ └─────────┘  │  ↓          ↓        │
                         │ RDD20     RDD27      │
┌─────────┐              │             ↓        │
│ RDD6    │              │           RDD31      │
│  ↓      │              │             ↓        │
│ RDD19   │──────────→   │           RDD40      │
└─────────┘              └──────────────────────┘
   S3
```

## Cost Lineage

*1. Keep stage information*

```
S1,2   RDD1
         ↓
S1,2   RDD2 ──────→  S2
         ↓          RDD11
S1     RDD5
```

# Collecting Potential Future References

## Submitted jobs (profiling phase)

S1
```
RDD1 → RDD2 → RDD5
```

S2
```
RDD1 → RDD2 → RDD11
```

S4
```
RDD1 → RDD2 → RDD17 → RDD20 → RDD21 → RDD27 → RDD31 → RDD40
```

S3
```
RDD6 → RDD19 → RDD20
```

## Cost Lineage

*1. Keep stage information*

S1,2    RDD1
S1,2    RDD2 → RDD11 (S2)
S1      RDD5
S3      RDD6
S3      RDD19

*2. Drop duplicate RDDs*

# Collecting Potential Future References

## Submitted jobs (profiling phase)

S1
```
RDD1
 ↓
RDD2
 ↓
RDD5
```

S2
```
RDD1
 ↓
RDD2
 ↓
RDD11
```

S4
```
RDD1
 ↓
RDD2
 ↓
RDD17
 ↓
RDD20 ──→ RDD21
              ↓
           RDD27
              ↓
           RDD31
              ↓
           RDD40
```

S3
```
RDD6
 ↓
RDD19 ──→ (RDD20)
```

## Cost Lineage

*1. Keep stage information*

***2. Drop duplicate RDDs***

```
S1,2,4   RDD1
           ↓
S1,2,4   RDD2 ──────────────→ RDD17   S4
           ↓        ↘            ↓
S1       RDD5      RDD11 S2    RDD20 ──→ RDD21   S4
           ↓               ↗    S4        ↓
S3       RDD6         ↗              RDD27   S4
           ↓     ↗                      ↓
S3       RDD19                       RDD31   S4
                                        ↓
                                     RDD40   S4
```

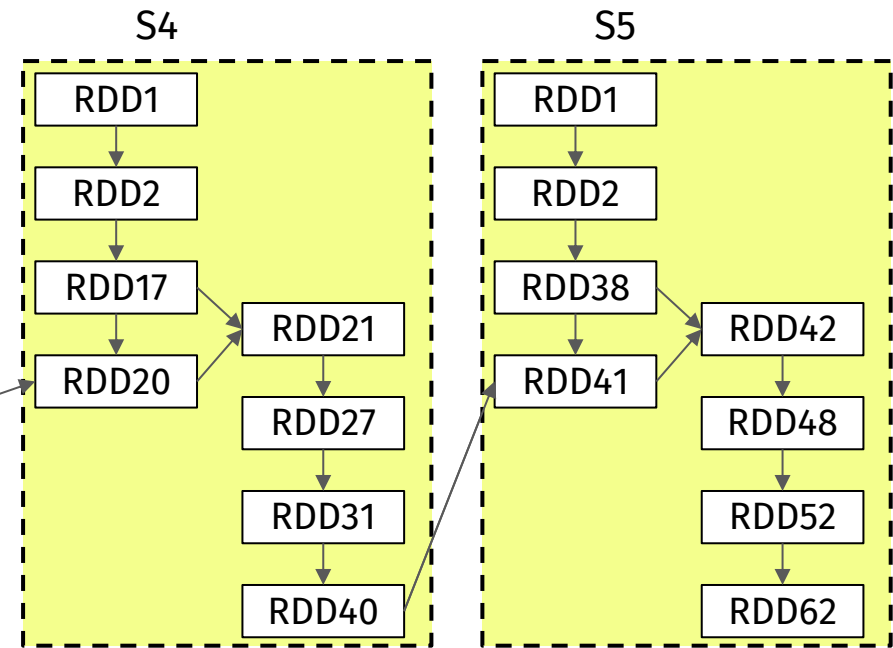# Collecting Potential Future References
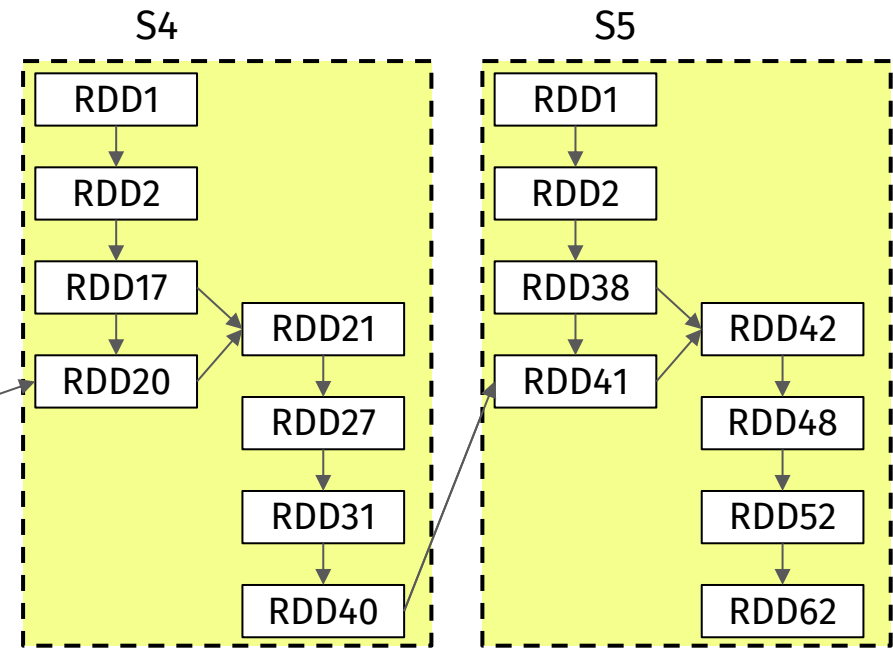
## Submitted jobs (profiling phase)

### S4

```
RDD1
  ↓
RDD2
  ↓
RDD17 ───→ RDD21
  ↓           ↓
RDD20 ────→ RDD27
            ↓
          RDD31
            ↓
          RDD40
```

### S5

```
RDD1
  ↓
RDD2
  ↓
RDD38 ───→ RDD42
  ↓           ↓
RDD41 ────→ RDD48
            ↓
          RDD52
            ↓
          RDD62
```

## Cost Lineage

*1. Keep stage information*

*2. Drop duplicate RDDs*

**3. Capture iteration loops**

S1,2,4,5 — RDD1

S1,2,4,5 — RDD2 → S4 RDD17

S1 — RDD5    S2 — RDD11 → RDD20 S4

S3 — RDD6

S3 — RDD19

S5 — RDD38

S5 — RDD41

RDD21 S4

RDD27 S4

RDD31 S4

RDD40 S4

# Collecting Potential Future References

## Submitted jobs

**S4**

RDD1 → RDD2 → RDD17 → RDD20
RDD17 → RDD21 → RDD27 → RDD31 → RDD40

**S5**

RDD1 → RDD2 → RDD38 → RDD41 → RDD42 → RDD48 → RDD52 → RDD62

## Cost Lineage (initial stages)

*1. Keep stage information*
*2. Drop duplicate RDDs*
*3. Capture iteration loops*

**4. Keep partition metrics**

S1,2,4,5

S1,2,4,5 — RDD2 —3s→ RDD17 (S4)

Time from **RDD2_partition1**: 3s

S1 RDD5   S2 RDD11   D20

**RDD17_Partition1** stored in executor1 deserialized memory, 100MB size

S3 RDD19

S5 RDD38 → RDD41
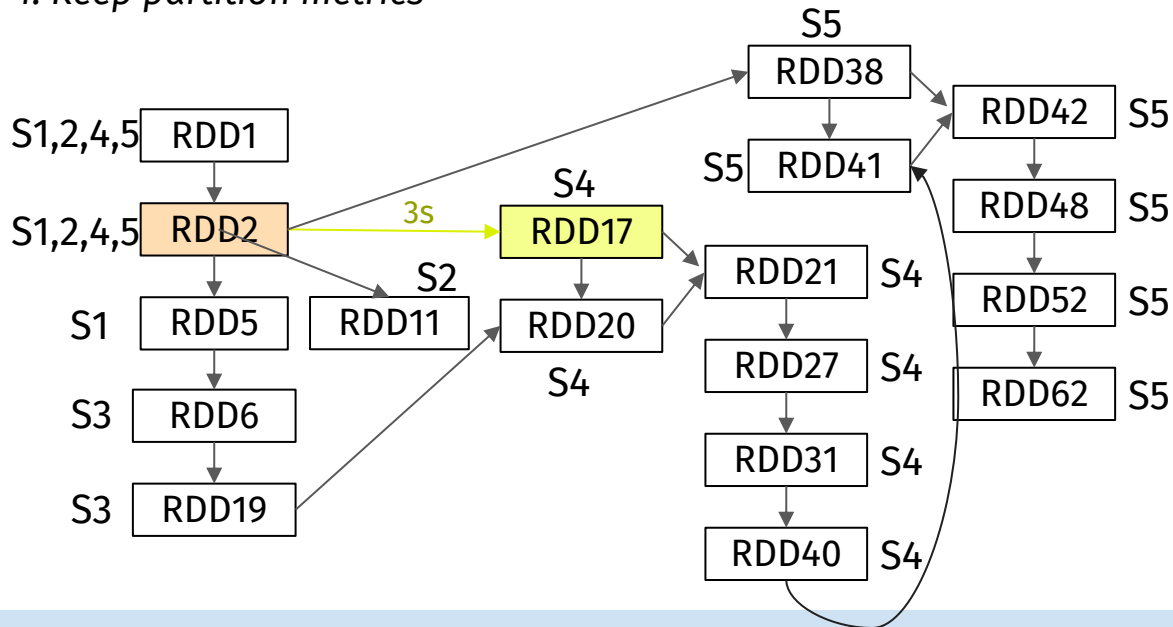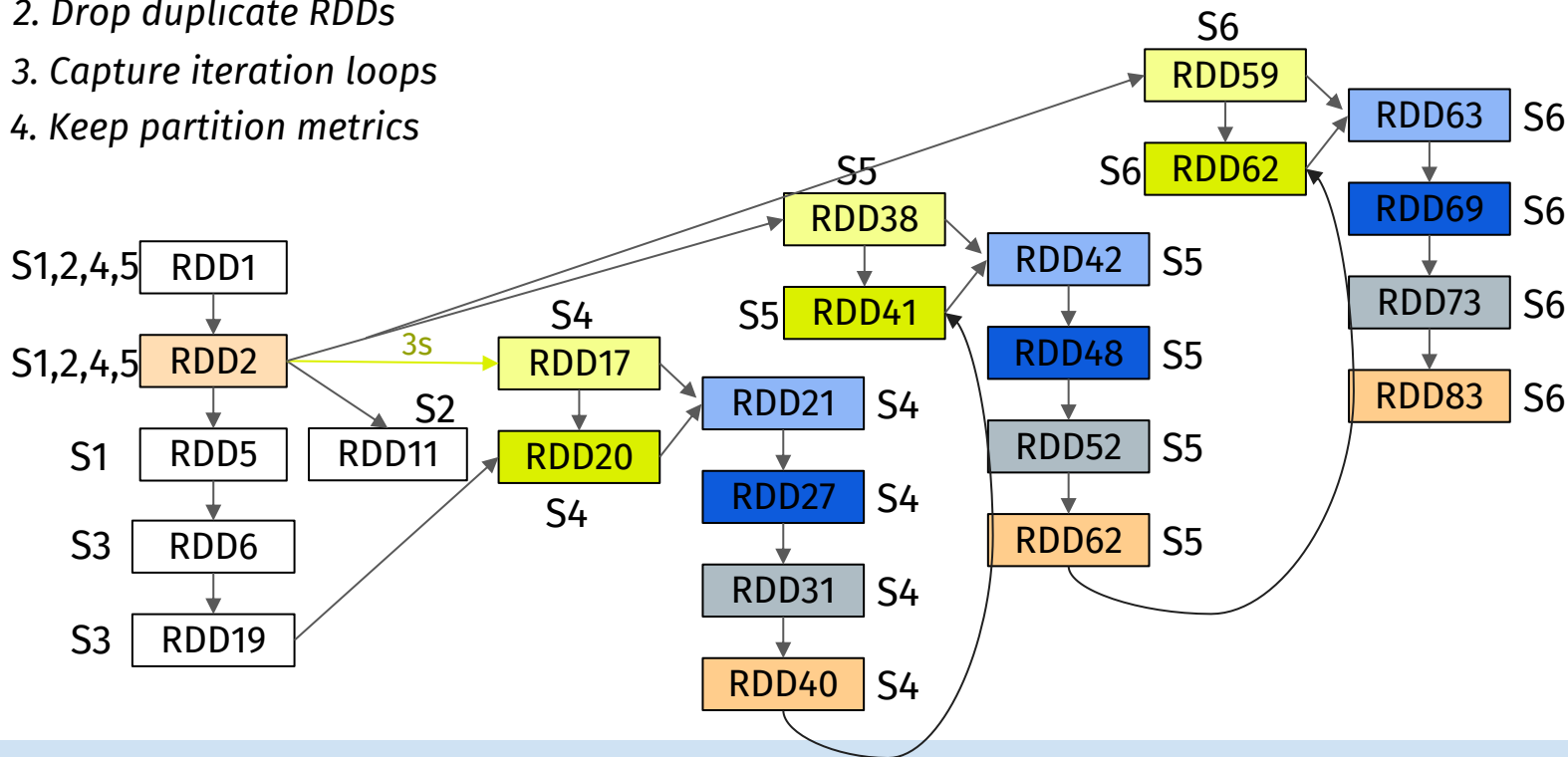RDD21 (S4) → RDD27 (S4) → RDD31 (S4) → RDD40 (S4)

# Collecting Potential Future References

## Cost Lineage

1. *Keep stage information*
2. *Drop duplicate RDDs*
3. *Capture iteration loops*
4. *Keep partition metrics*

**5. *Perform induction on future iterations***

# Collecting Potential Future References

**Cost Lineage**

*1. Keep stage information*

*2. Drop duplicate RDDs*

*3. Capture iteration loops*

*4. Keep partition metrics*
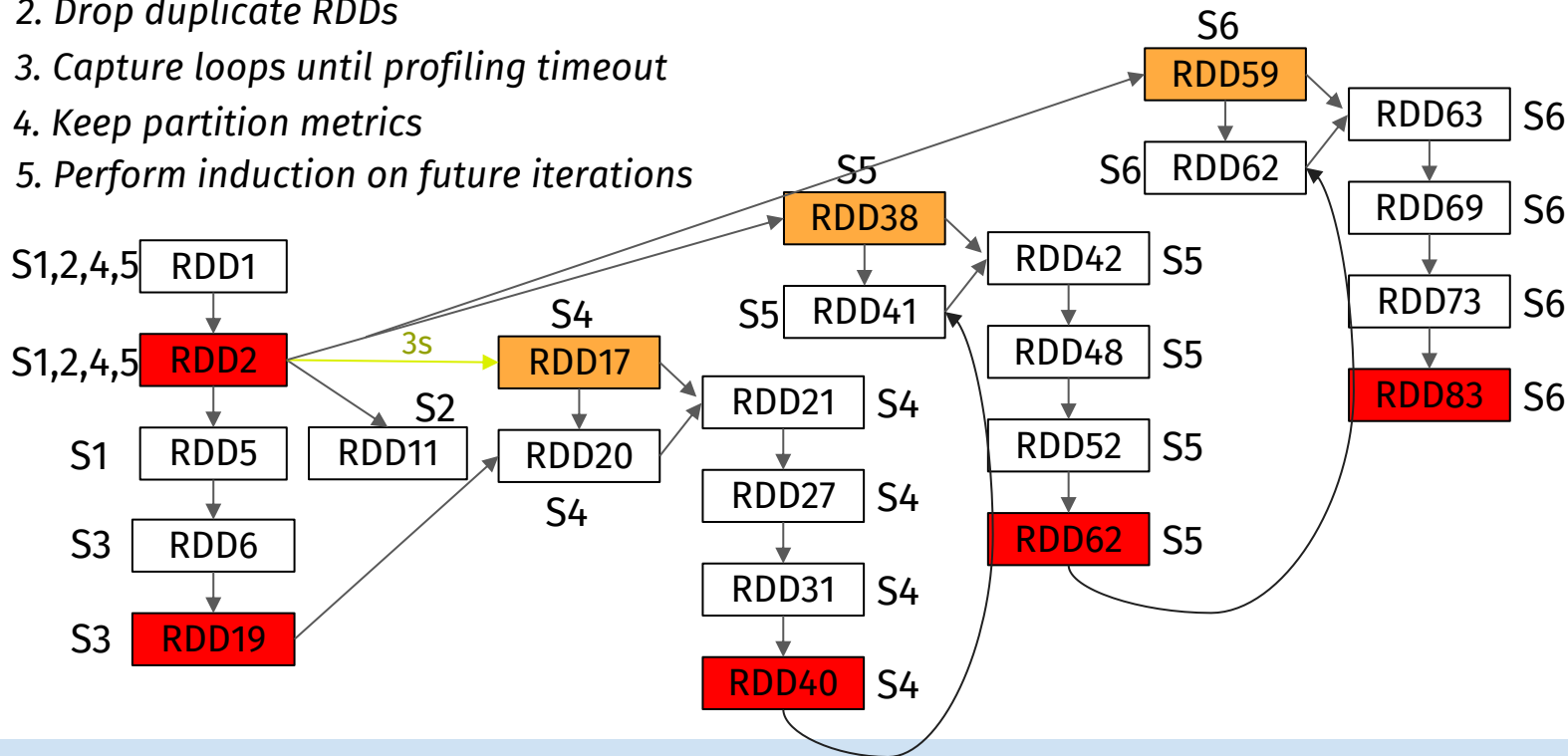
**5. Perform induction on future iterations**

# Collecting Potential Future References

## Cost Lineage

1. *Keep stage information*
2. *Drop duplicate RDDs*
3. *Capture loops until profiling timeout*
4. *Keep partition metrics*
5. *Perform induction on future iterations*

**6. Automatic caching decisions
(Triggered after each stage execution)**
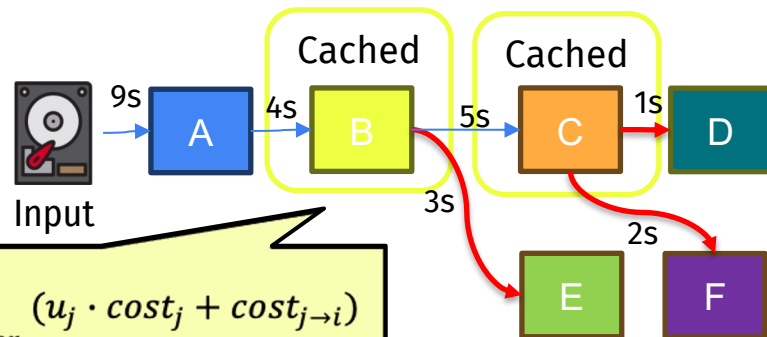
# Potential Cost Estimation

**Disk Cost**

$$cost_{i_d} = \frac{size_i}{throughput_{disk}}$$

| | 5 ∨ | 1000MB ∨ | C: 37% (83/226GB) ∨ |
|---|---|---|---|
| | | Read [MB/s] | Write [MB/s] |
| Seq | | 458.4 | 256.5 |
| 512K | | 380.7 | 245.0 |
| 4K | | 22.77 | 97.87 |
| 4K QD32 | | 194.2 | 225.6 |

**Recomputation Cost**



Input

Cached   Cached

9s → A → 4s → B → 5s → C → 1s → D

3s → E      2s → F

$$cost_{i_r} = \max_{j \in P_{ancestor_i}} (u_j \cdot cost_j + cost_{j \to i})$$

Linearly increases with the data size
(Depends on disk performance, like RPM)

Depends on future data dependencies

Depends on cached ancestor data,
which varies at runtime

# Potential Cost Estimation

**Disk Cost**

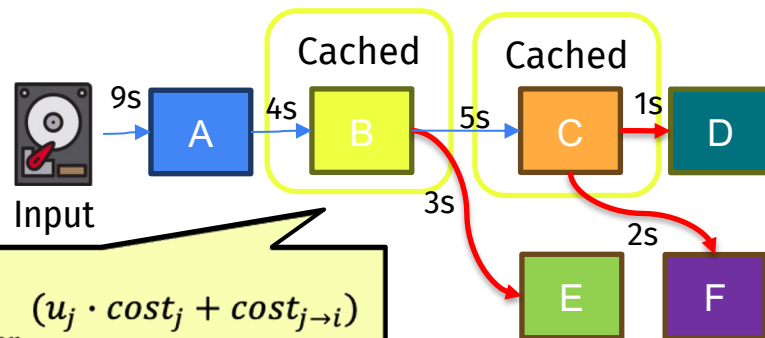$$cost_{i_d} = \frac{size_i}{throughput_{disk}}$$

| All | 5 | 1000MB | C: 37% (83/226GB) |
|-----|---|--------|-------------------|
| | Read [MB/s] | | Write [MB/s] |
| Seq | 458.4 | | 256.5 |
| 512K | 380.7 | | 245.0 |
| 4K | 22.77 | | 97.87 |
| 4K QD32 | 194.2 | | 225.6 |

**Recomputation Cost**



$$cost_{i_r} = \max_{j \in P_{ancestor_i}} (u_j \cdot cost_j + cost_{j \to i})$$

Linearly increases with the data size
(Depends on disk performance, like RPM)

Depends on future data dependencies

Depends on cached ancestor data,
which varies at runtime

# Potential Cost Estimation

**Disk Cost**

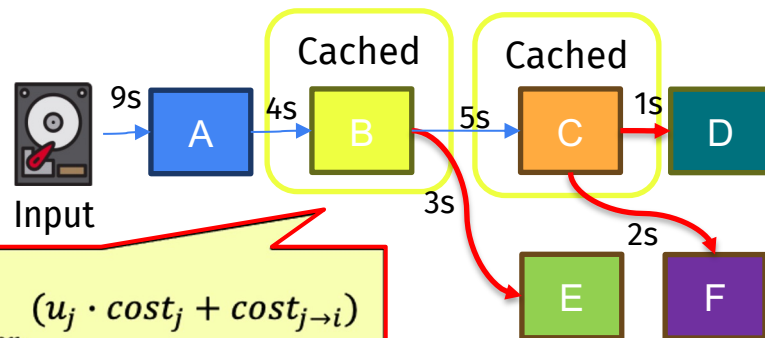$$cost_{i_d} = \frac{size_i}{throughput_{disk}}$$

| All | 5 ⌄ | 1000MB ⌄ | C: 37% (83/226GB) ⌄ | |
|---|---|---|---|---|
| | | **Read [MB/s]** | **Write [MB/s]** | |
| Seq | | 458.4 | 256.5 | |
| 512K | | 380.7 | 245.0 | |
| 4K | | 22.77 | 97.87 | |
| 4K QD32 | | 194.2 | 225.6 | |

Linearly increases with the data size
(Depends on disk performance, like RPM)

**Recomputation Cost**

Cached    Cached

Input    9s → A → 4s → B → 5s → C → 1s → D

3s    E    2s    F
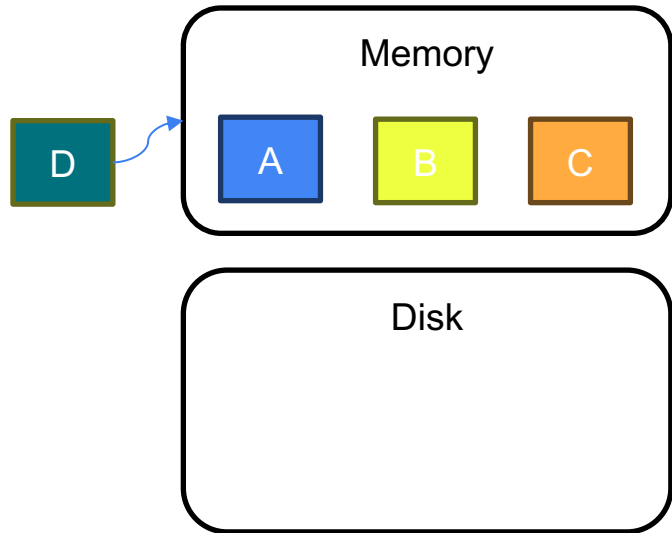
$$cost_{i_r} = \max_{j \in P_{ancestor_i}} (u_j \cdot cost_j + cost_{j \to i})$$

Depends on future data dependencies

Depends on cached ancestor data,
which varies at runtime

# Decision Making Algorithm (Trigger)

Algorithm triggered after each iteration when:

1) size of new data to be cached exceeds free memory
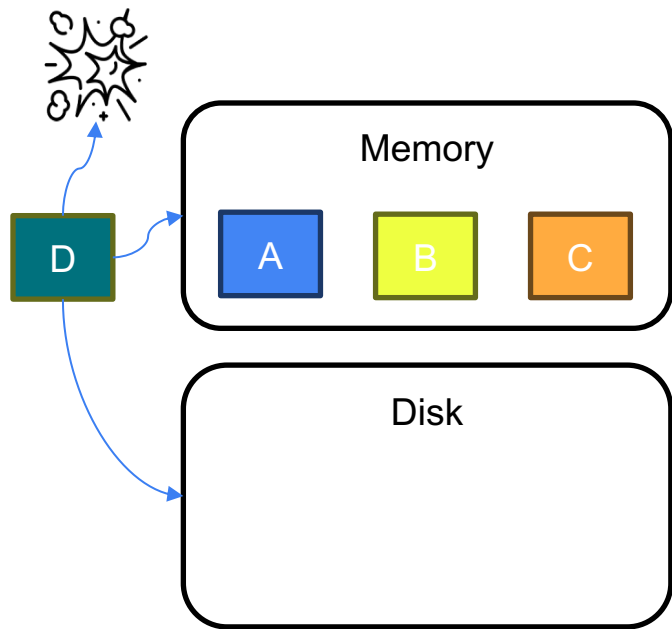
$$size_{new} > free\_capacity_{mem}$$

AND

2) expected potential recovery overhead of new data to be prevented (min of recomputation cost and disk cost)

$$cost_{new} = \min(cost_{new_r},\ cost_{new_d})$$

is larger than the potential overhead of any data in memory:

$$cost_{new} > \min_{d\,\in\,D_{mem}} cost_d$$

# Decision Making Algorithm (ILP)



ILP Solver

Action space: mem-cached / disk-cached / un-cached state, for each partition

$$\sum_{i \in P} m_i + \sum_{i \in P} d_i + \sum_{i \in P} u_i = |P| \ \text{AND} \ \forall i \in P, \ m_i + d_i + u_i = 1$$
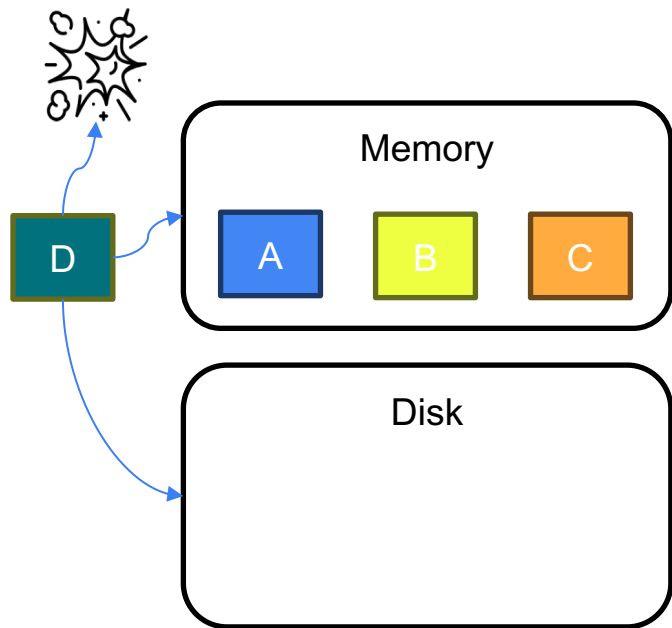
Constraint: memory space does not exceed sum of mem-cached partition data

$$\sum_{i \in P} size_i \cdot m_i \leq capacity_{mem}$$

Objective: minimize the sum of potential recovery overhead and JCT for the RDDs in the current & next job $J$ (upcoming iterations)

$$\text{argmin} \sum_{i \in J} cost_i = \text{argmin} \sum_{i \in J} (u_i \cdot cost_{i_r} + d_i \cdot cost_{i_d})$$

# Decision Making Algorithm

ILP Solver

Action space: mem-cached / disk-cached / un-cached state, for each partition

$$\sum_{i \in P} m_i + \sum_{i \in P} d_i + \sum_{i \in P} u_i = |P| \text{ AND } \forall i \in P,\ m_i + d_i + u_i = 1$$

Constraint: memory space does not exceed sum of mem-cached partition data

$$cost_{i_r} = \max_{j \in P_{ancestor_i}} (u_j \cdot cost_j + cost_{j \to i})$$

$$cost_{i_d} = \frac{size_i}{throughput_{disk}}$$

Objective: minimize the sequential recovery overhead & JCT for the RDDs in the current & next ~~batch J~~ (upcoming iteration)

$$\text{argmin} \sum_{i \in J} cost_i = \text{argmin} \sum_{i \in J} (u_i \cdot cost_{i_r} + d_i \cdot cost_{i_d})$$

**Memory**

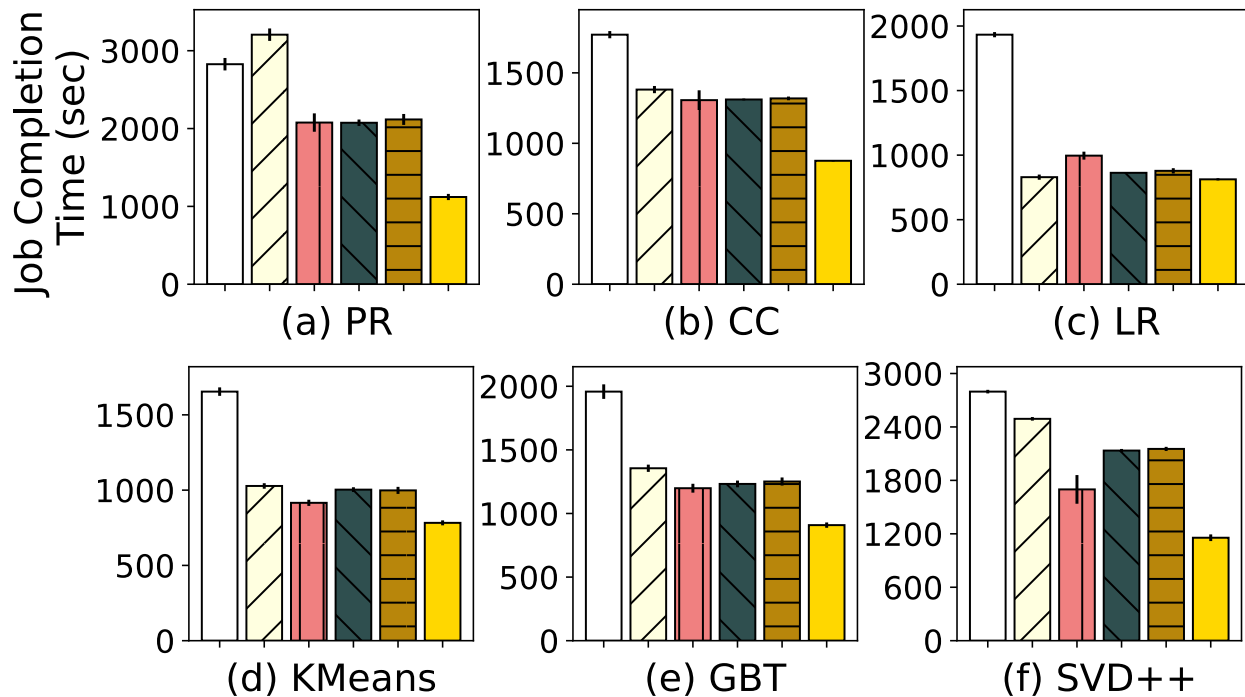A   B   C

**Disk**

D

# Blaze Implementation

- Data processing runtime

  - Spark 3.3.2, modified and added 6K lines of Scala 2.12 code

- Blaze profiler & cache optimizer

  - ~500 lines of Bash script
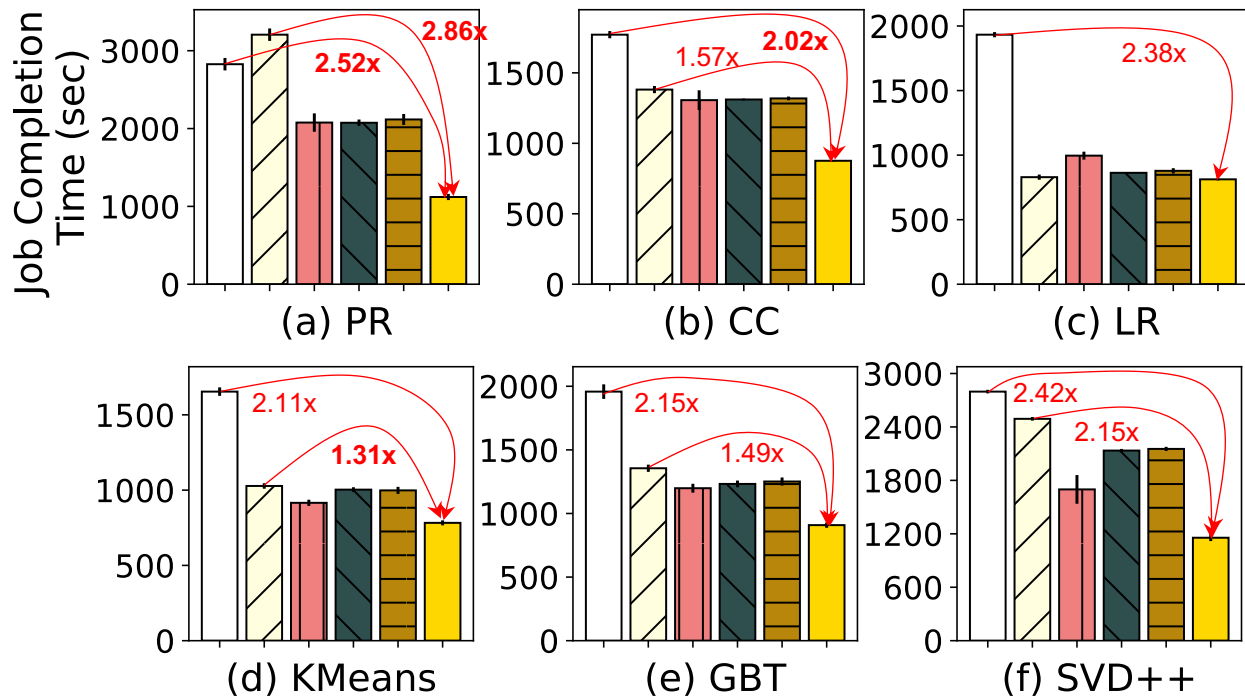
- ILP solver

  - Gurobi optimizer 10.0.1

# Evaluation Setup

- 11 instances (1 as master, 10 as executors)
  - AWS EC2 r5.2xlarge instances (8vCores, 64GB memory, 10Gbps network, each)
    - → A total of **80** executor vCores, **500GB** executor memory (**170GB** memory used as **cache memory**).
  - **100GB SSD** (gp2) on each instance for caching stores (disks) → A total of **1TB SSD** for a **disk** caching store

- Workloads
  - **PageRank (PR)**
  - **Connected Components (CC)**
  - **Logistic Regression (LR)**
  - **K-Means Clustering (KMeans)**
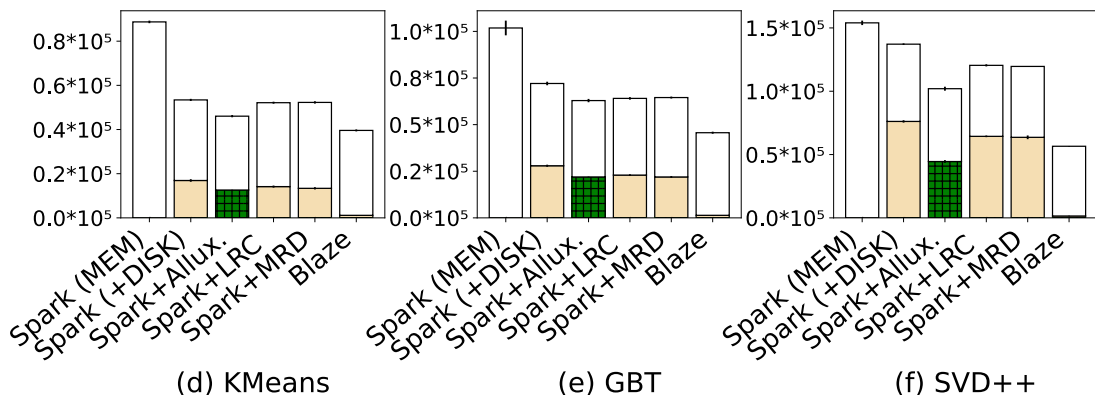  - **Gradient Boosted Trees (GBT)**
  - **Singular Vector Decomposition++ (SVD++)**

# End-to-End Performance Comparison



Legend: Spark (MEM), Spark (MEM+DISK), Spark+Alluxio, LRC, MRD, Blaze

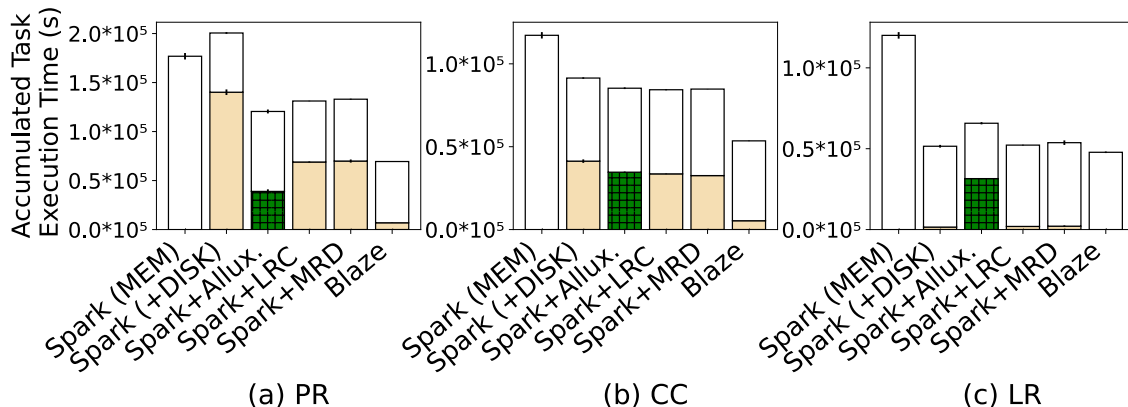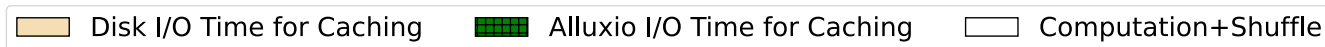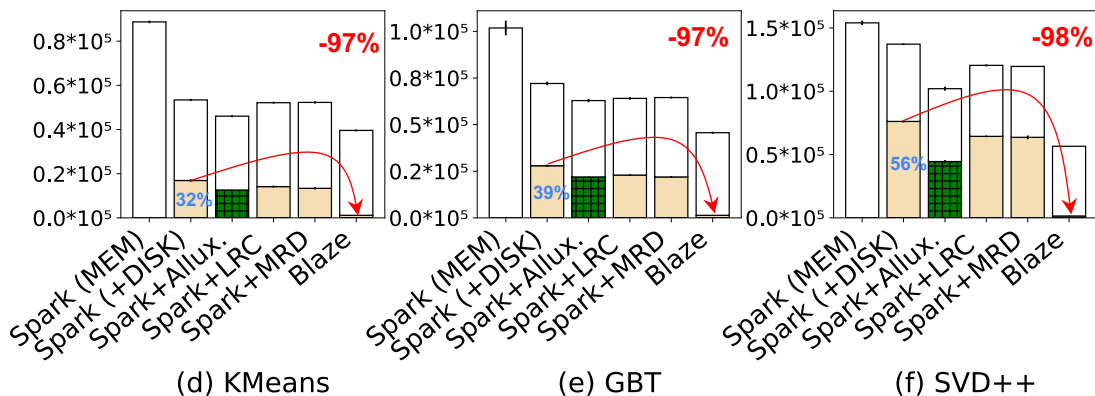(a) PR  (b) CC  (c) LR  (d) KMeans  (e) GBT  (f) SVD++

Job Completion Time (sec)

# End-to-End Performance Comparison



Legend: Spark (MEM), Spark (MEM+DISK), Spark+Alluxio, LRC, MRD, Blaze

(a) PR — 2.52x, 2.86x
(b) CC — 1.57x, 2.02x
(c) LR — 2.38x
(d) KMeans — 2.11x, 1.31x
(e) GBT — 2.15x, 1.49x
(f) SVD++ — 2.42x, 2.15x

Y-axis: Job Completion Time (sec)

# Performance Analysis: Computation+Shuffle vs. Disk I/O



Legend: Disk I/O Time for Caching | Alluxio I/O Time for Caching | Computation+Shuffle

(a) PR     (b) CC     (c) LR

(d) KMeans     (e) GBT     (f) SVD++

# Performance Analysis: Computation+Shuffle vs. Disk I/O



Legend: Disk I/O Time for Caching | Alluxio I/O Time for Caching | Computation+Shuffle

(a) PR  (b) CC  (c) LR  (d) KMeans  (e) GBT  (f) SVD++

# Blaze: Performance Breakdown



Legend: Spark (MEM+DISK), +AutoCache, +CostAware, Blaze

(a) PR  (b) CC  (c) LR

(d) KMeans  (e) GBT  (f) SVD++

# Blaze: Performance Breakdown



(a) PR  (b) CC  (c) LR

(d) KMeans  (e) GBT  (f) SVD++

# Blaze Summary

- Caching is crucial in reducing the recomputation costs for iterative data processing workloads (e.g., graph processing, ML)

- Existing separated, greedy mechanisms lead to unnecessary caching and inefficient use of memory space

- Recomputation and disk overheads + potential references have to be tracked dynamically

- Based on the tracked information, Blaze makes sophisticated decisions on which data to cache and to evict with an ILP solver

- Blaze achieves up to 2.86x speedup on end-to-end performance and optimizes cache data by 95% on average

# Thank you!

## Blaze:
## Holistic Caching for Iterative Data Processing



**Contact me for further questions**

wonook@apache.org                              https://wonook.github.io