

# ***TraceTracker:*** Hardware/Software Co-Evaluation for Large-Scale I/O Workload Reconstruction

**Miryeong Kwon, Jie Zhang, Gyuyoung Park,  
Wonil Choi, David Donofrio, John Shalf,  
Mahmut Kandemir, and Myoungsoo Jung**



**CAMEL**ab.org  
Computer Architecture and  
Memory Systems  
Laboratory

# Summary

---

- **Motivation:** Block traces collected on old systems (a decade ago HDDs) can make system analysis significantly different from the modern systems
- **Challenges:** Trace reconstruction is high-cost work and inaccurate
  - Application execution requires prohibitive resources; thousands of users, large-scale computing systems (**impractical**)
  - Overly-simplified methods are **excessively imprecise** due to lack of runtime contexts
- **Goal:** **Accurately** reconstruct workloads **without** application execution
- **Solutions:** Hardware/Software co-evaluation method
  1. Software method: Infer runtime contexts from existing block traces
  2. Hardware method: Remaster storage traces; aware of target system
- **Evaluation:**
  - TraceTracker infer runtime contexts 99%, 96% for number of occurrences and total periods, respectively
  - Apply TraceTracker to 577 open-license block traces



# TraceTracker

---

**1. Background of Block Trace**

**2. Trace Reconstruction Methods**

**3. Insights on Trace Reconstruction**

**4. TraceTracker Method**

**5. Evaluation**



# Block Trace

---

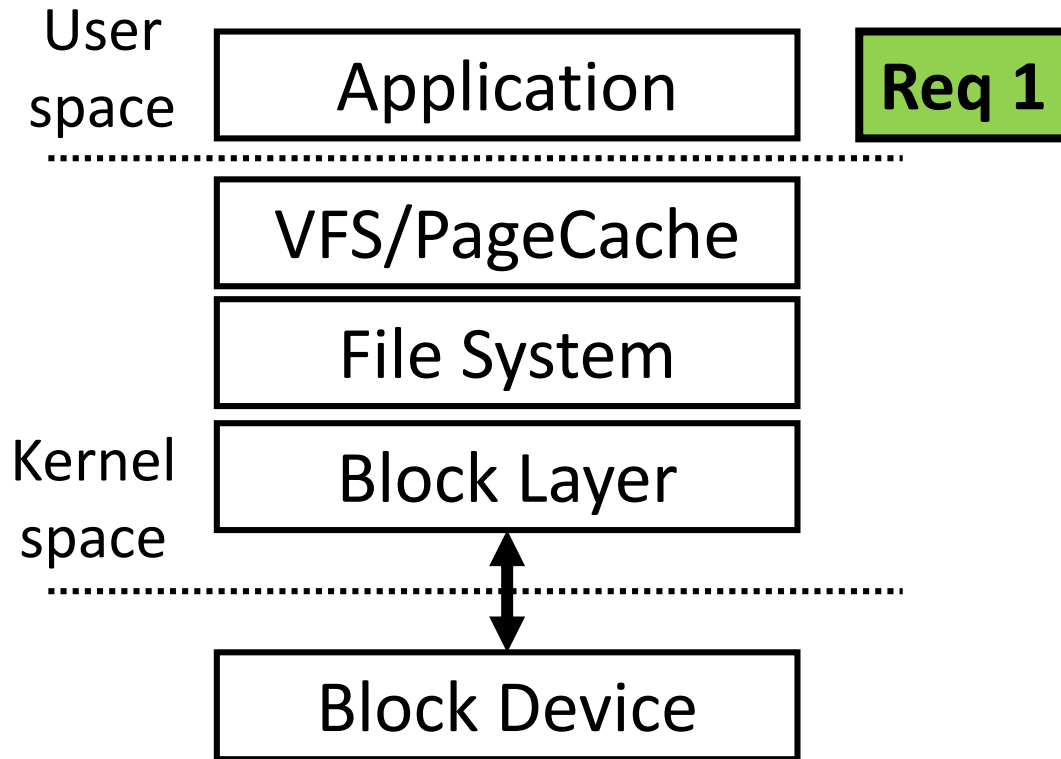
: Consists of detailed **block device access information**



# Block Trace

---

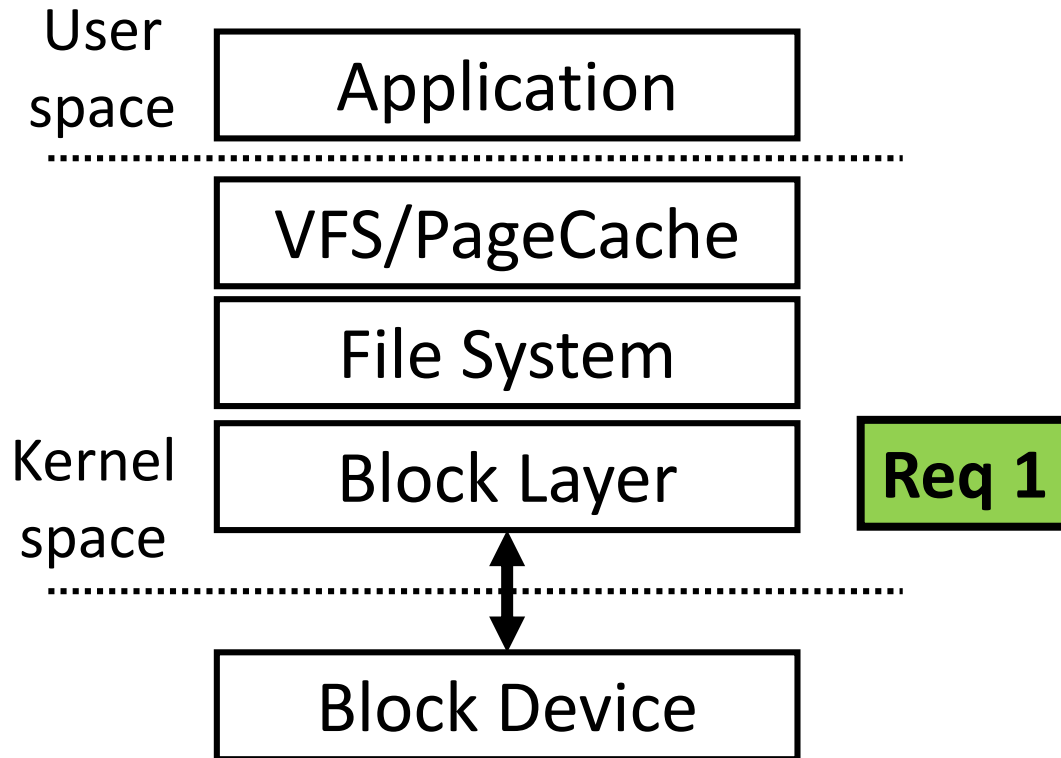
: Consists of detailed **block device access information**



# Block Trace

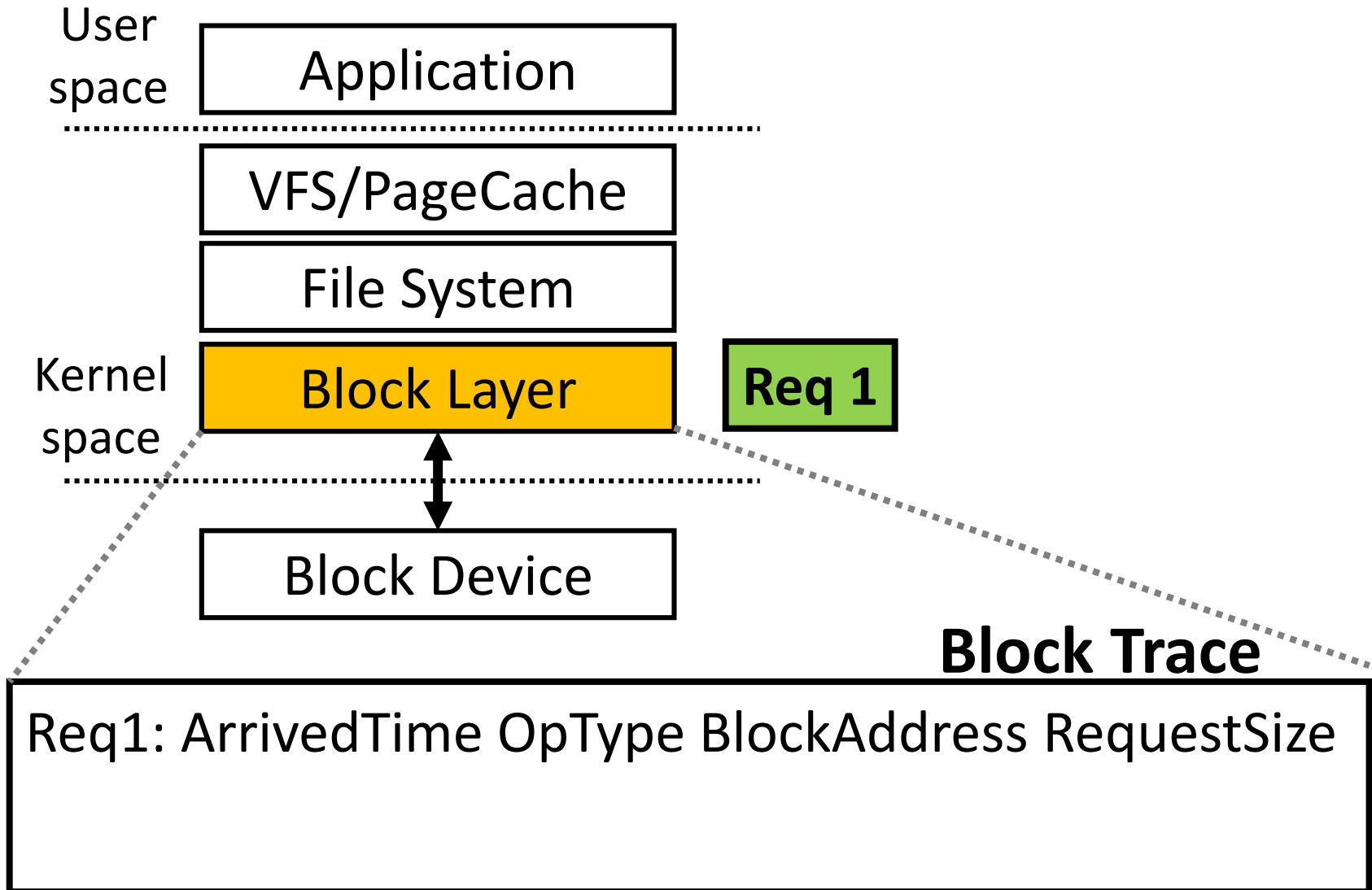
---

: Consists of detailed **block device access information**



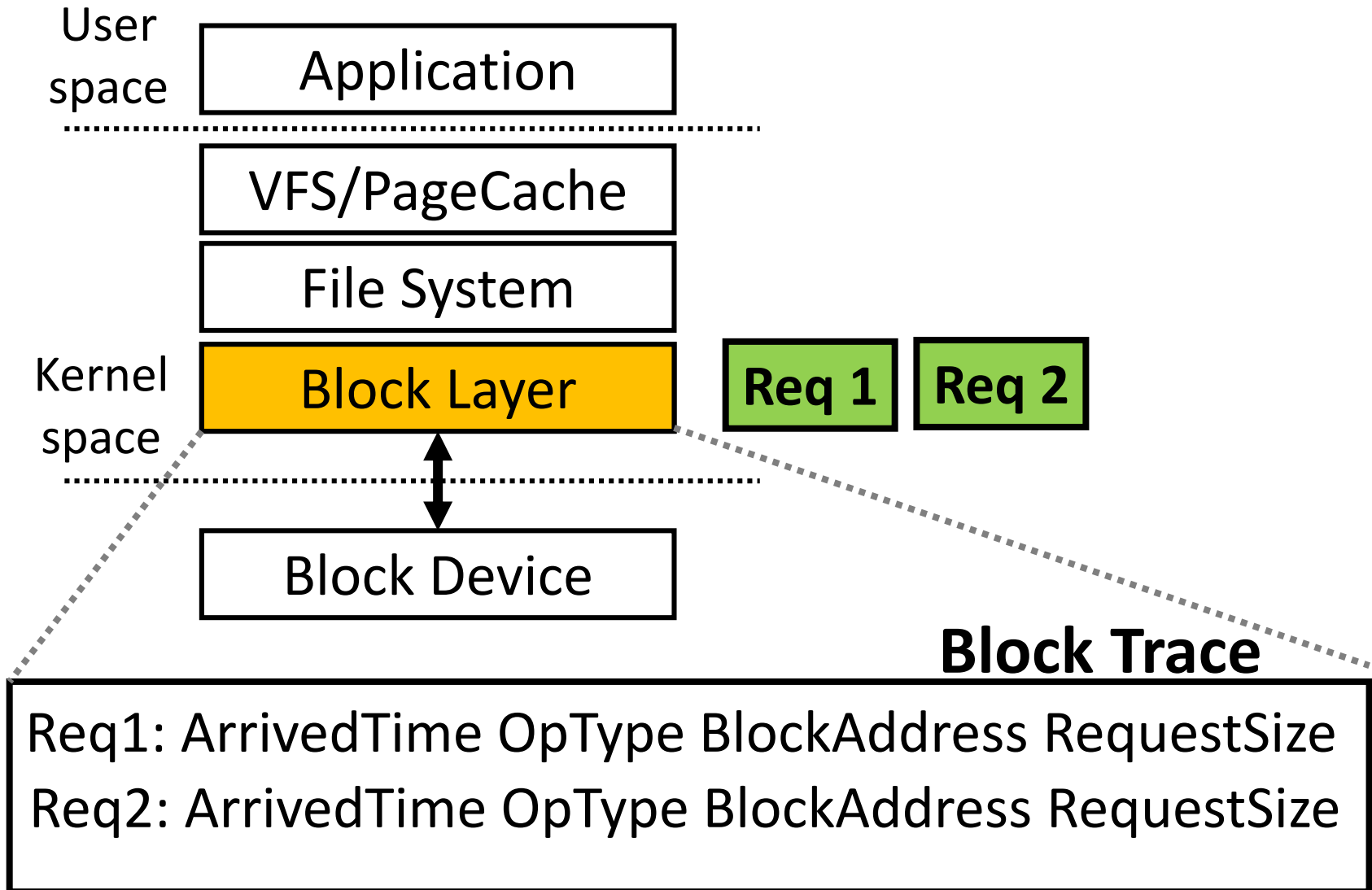
# Block Trace

: Consists of detailed **block device access information**



# Block Trace

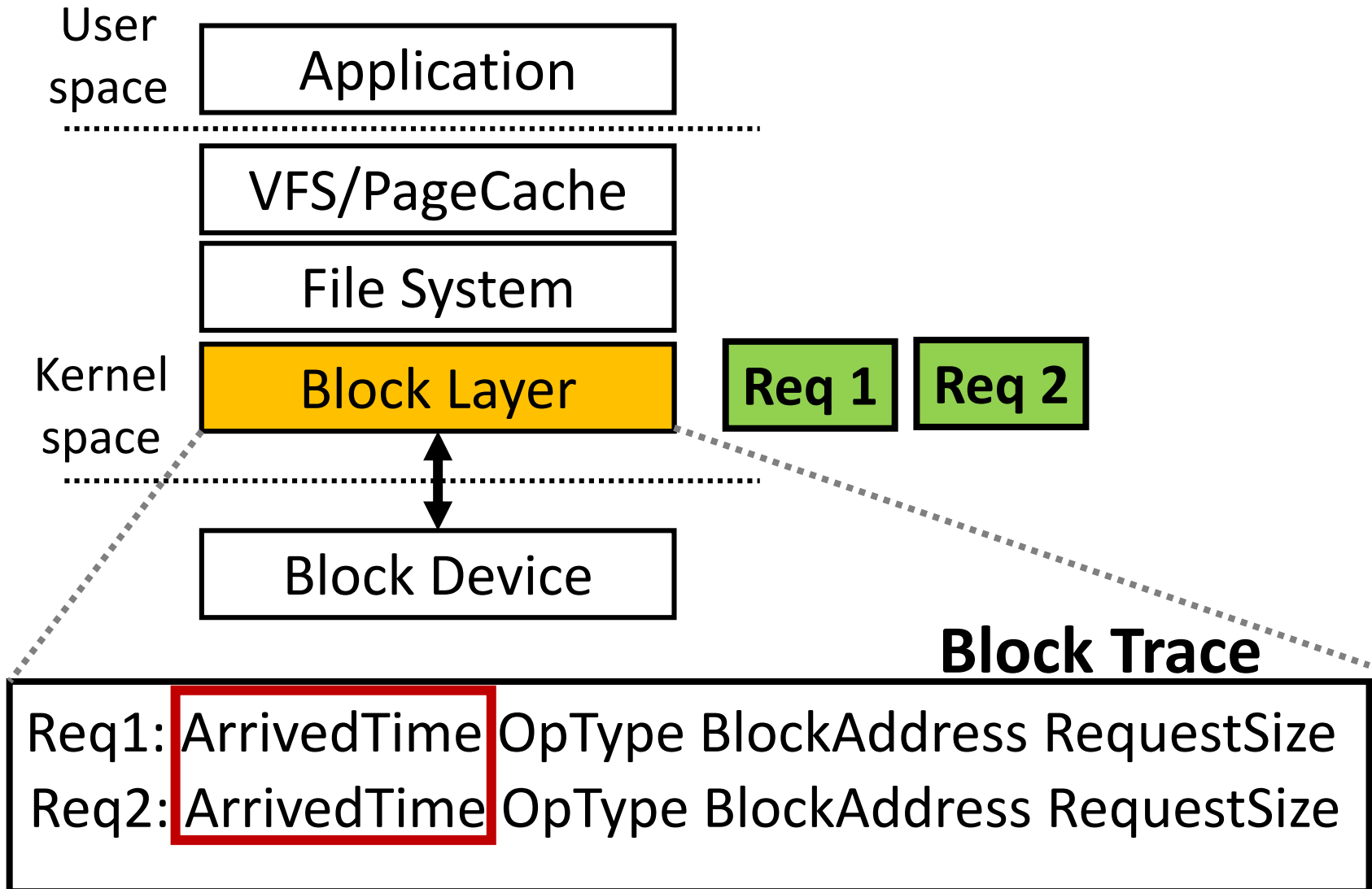
: Consists of detailed **block device access information**





# Block Trace

: Consists of detailed **block device access information**



# Inaccuracy of existing block traces

---

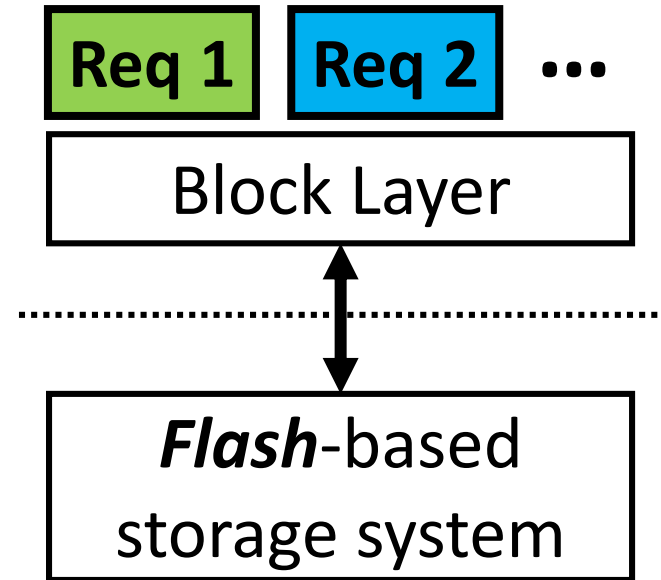
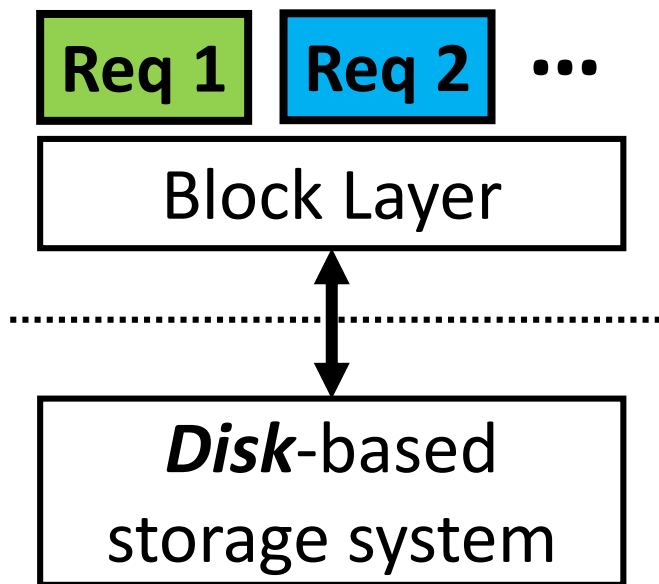
**Existing traces:** collected on the *disk-based* storage systems and published around *2007*



# Inaccuracy of existing block traces

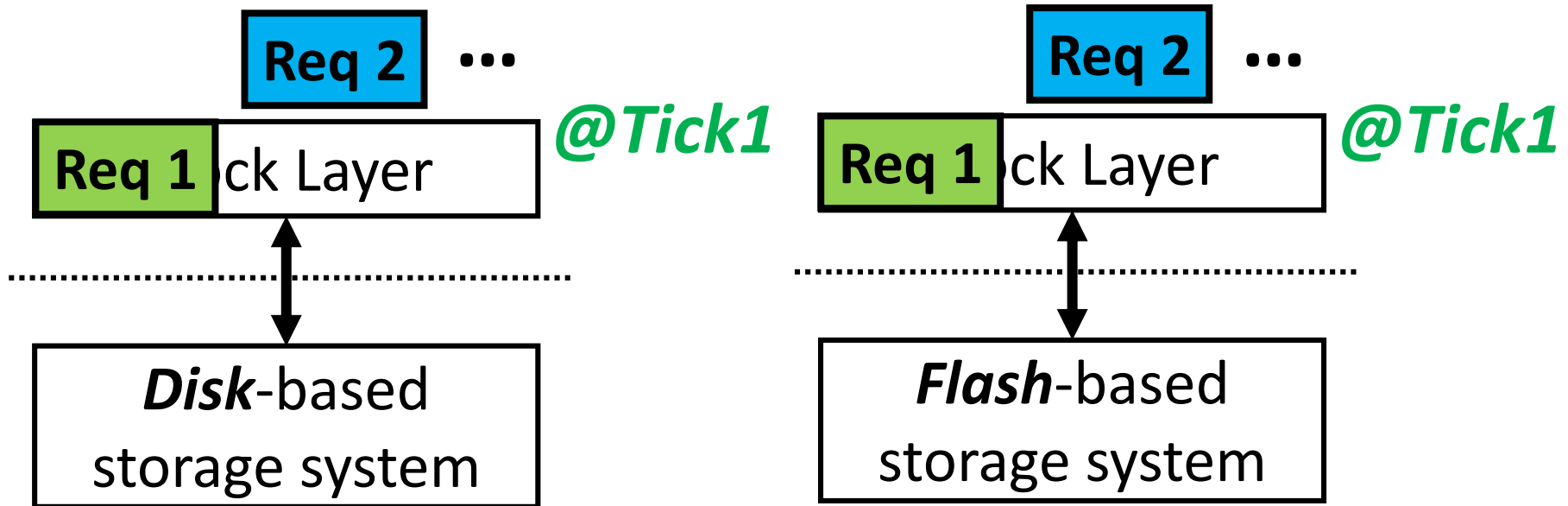
---

**Existing traces:** collected on the *disk-based* storage systems and published around *2007*



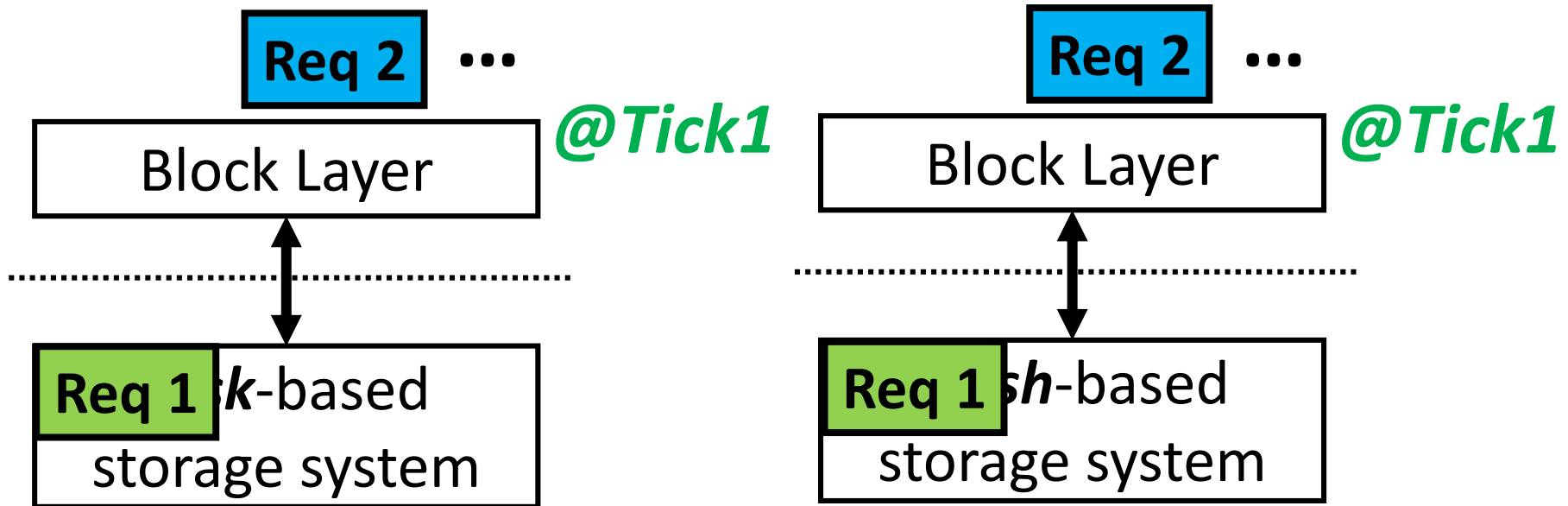
# Inaccuracy of existing block traces

**Existing traces:** collected on the *disk-based* storage systems and published around *2007*



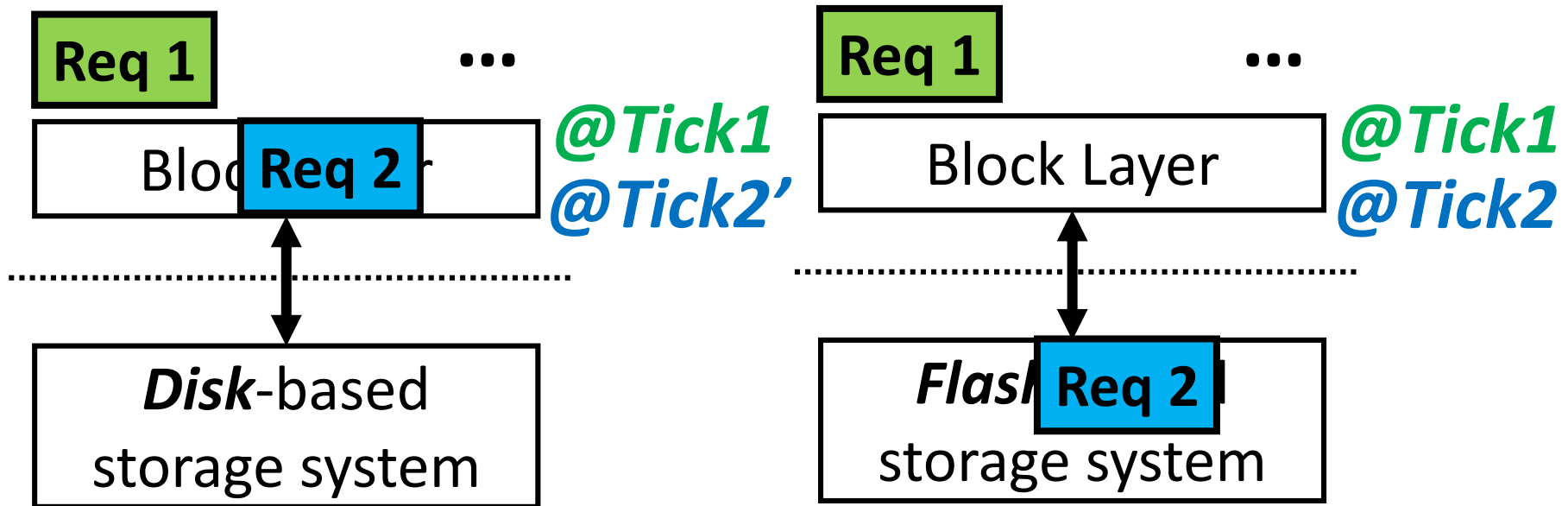
# Inaccuracy of existing block traces

**Existing traces:** collected on the *disk-based* storage systems and published around *2007*



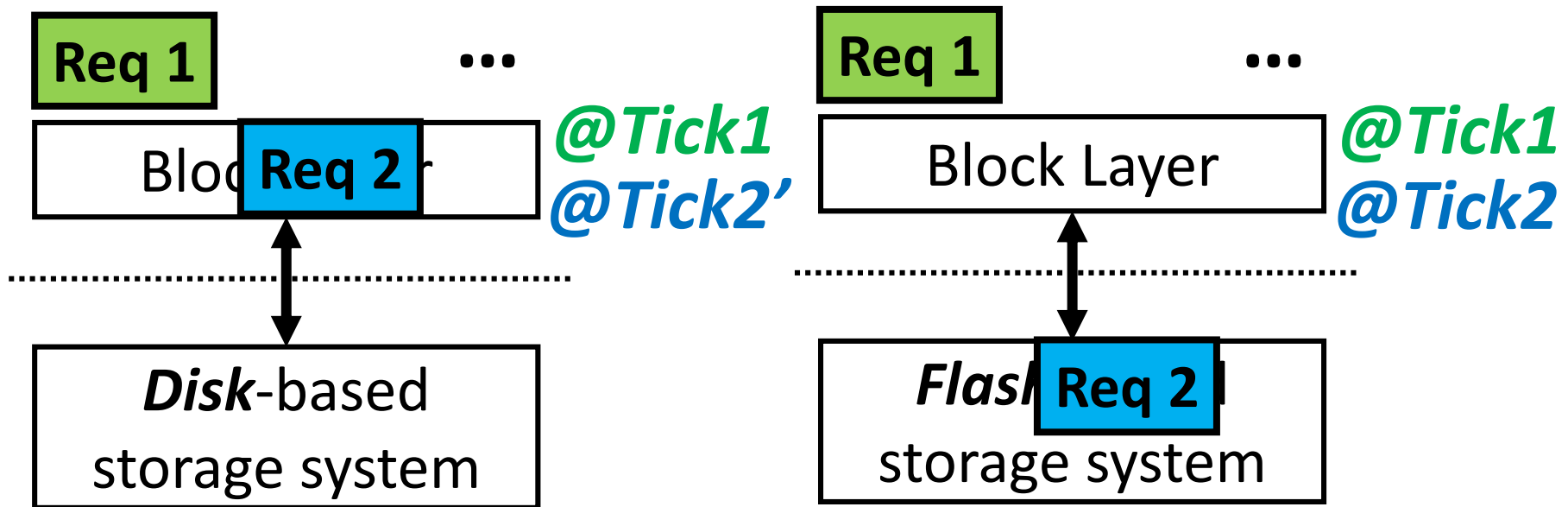
# Inaccuracy of existing block traces

**Existing traces:** collected on the *disk-based* storage systems and published around *2007*



# Inaccuracy of existing block traces

**Existing traces:** collected on the *disk-based* storage systems and published around *2007*



Timing information of block traces should be revised correspond to target storage system



# TraceTracker

---

1. Background of Block Trace

**2. Trace Reconstruction Methods**

3. Insights on Trace Reconstruction

4. TraceTracker Method

5. Evaluation





# Trace Reconstruction Methods

---



# Trace Reconstruction Methods

---

- *APPLICATION*

Large-scale  
computing  
system

Target  
Storage System



# Trace Reconstruction Methods

---

- ***APPLICATION***

Target  
Application

Large-scale  
computing  
system

Target  
Storage System



# Trace Reconstruction Methods

---

## • *APPLICATION*

Target  
Application

Large-scale  
computing  
system

Target  
Storage System

## • *REVISION*

Old Block Trace

Target  
Storage System



# Trace Reconstruction Methods

---

## • *APPLICATION*

Target  
Application

Large-scale  
computing  
system

Target  
Storage System

## • *REVISION*

Old Block Trace

Target  
Storage System

## • *ACCELERATION*

Old Block Trace

$$\text{Acceleration ratio} = \frac{\text{Bandwidth}_{old}}{\text{Bandwidth}_{new}}$$



# Trace Reconstruction Methods

## • APPLICATION

Target  
Application

Large-scale  
computing  
system

Target  
Storage System

## • REVISION

Old Block Trace

Target  
Storage System

## • ACCELERATION

Old Block Trace

$$\text{Acceleration ratio} = \frac{\text{Bandwidth}_{old}}{\text{Bandwidth}_{new}}$$

Accuracy:

>>

>>

Cost:

>>>

>>



# Trace Reconstruction Methods

## • APPLICATION

Target  
Application

Large-scale  
computing  
system

Target  
Storage System

## • REVISION

Old Block Trace

Do not include runtime context

Target  
Storage System

## • ACCELERATION

Old Block Trace

$$\frac{Bandwidth_{old}}{Bandwidth_{new}}$$

Accuracy:

>>

>>

Cost:

>>>

>>



# Trace Reconstruction Methods

## • APPLICATION

## • REVISION

## • ACCELERATION

Target  
Application



*User idle times*

Large-scale  
computing  
system

Target  
Storage System

Old Block Trace

Old Block Trace

**Do not include runtime context**

Target  
Storage System

$$\frac{Bandwidth_{old}}{Bandwidth_{new}}$$

Accuracy:

>>

>>

Cost:

>>>

>>





# Trace Reconstruction Methods

## • APPLICATION

## • REVISION

## • ACCELERATION

Target  
Application

*User idle times*

Large-scale  
computing  
system

*System delay times*

Target  
Storage System

Old Block Trace

Old Block Trace

**Do not include runtime context**

Target  
Storage System

$$\frac{Bandwidth_{old}}{Bandwidth_{new}}$$

Accuracy:

>>

>>

Cost:

>>>

>>



# Trace Reconstruction Methods

## • APPLICATION

Target  
Application

Large-scale  
computing  
system

Target  
Storage System

## • REVISION

Old Block Trace

Target  
Storage System

## • ACCELERATION

Old Block Trace

Lack of detailed  
timing  
information  
 $\text{Acceleration ratio} = \frac{\text{Bandwidth}_{\text{old}}}{\text{Bandwidth}_{\text{new}}}$

Accuracy:

>>

>>

Cost:

>>>

>>

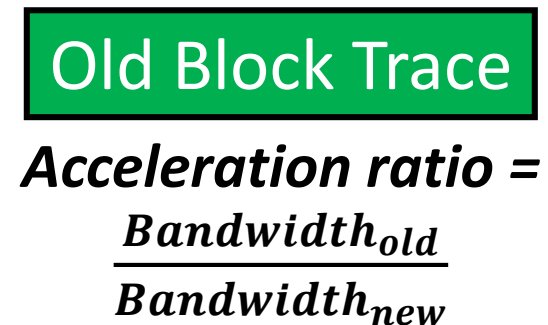
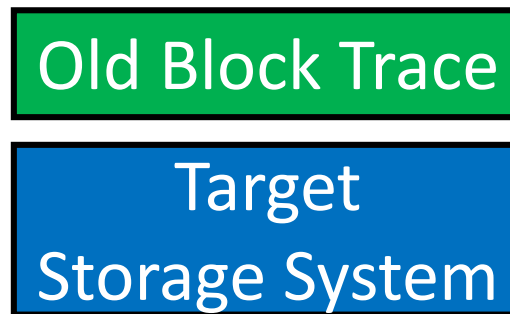
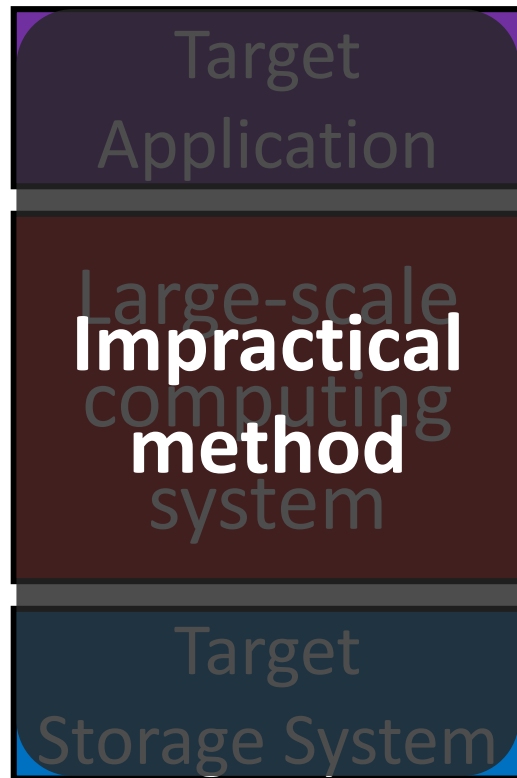


# Trace Reconstruction Methods

• **APPLICATION**

• **REVISION**

• **ACCELERATION**



Accuracy:

>>

>>

Cost:

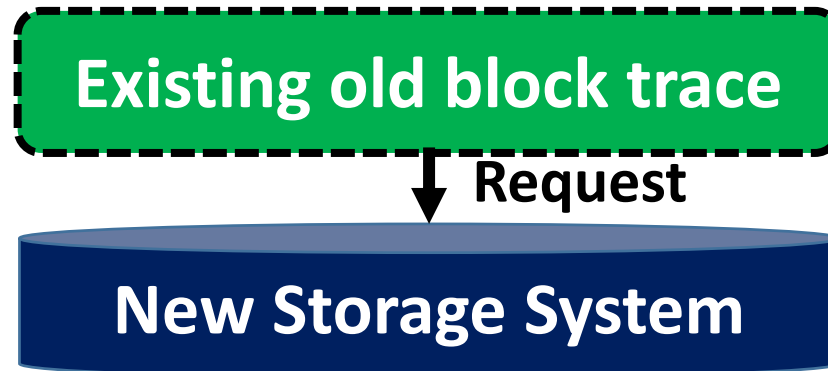
>>>

>>



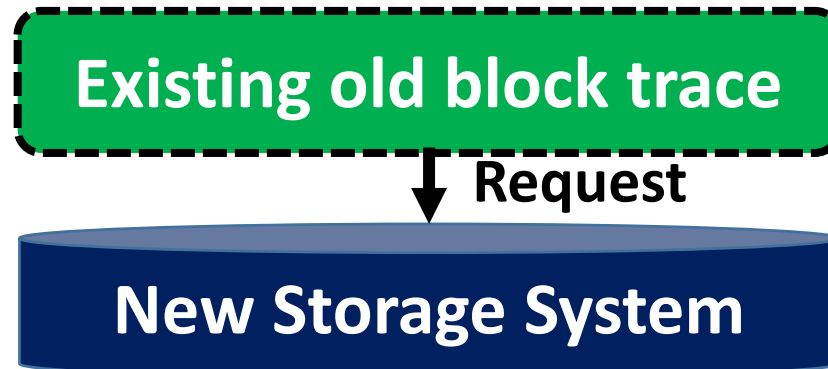
# Goal of TraceTracker

---



# Goal of TraceTracker

---

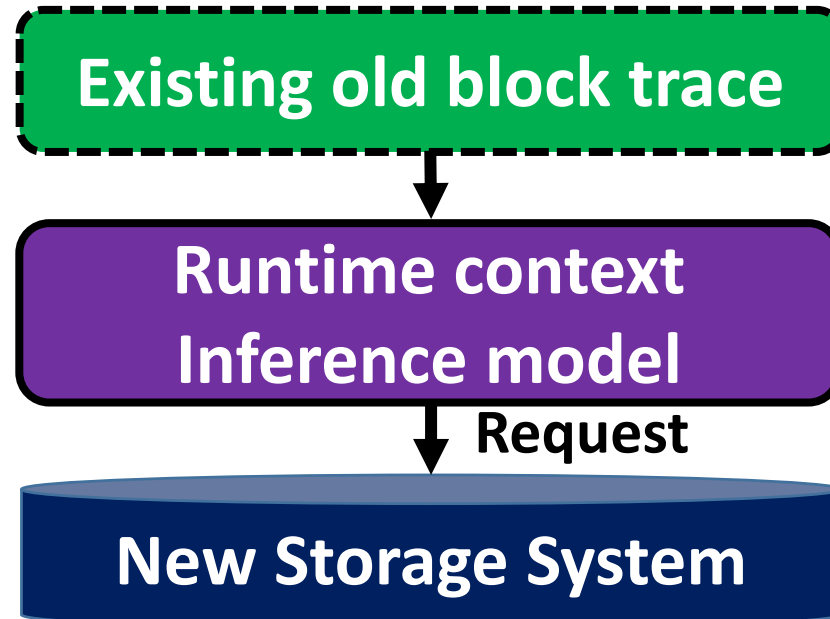


**Goal: Accurately reconstruct workloads  
without application execution (at low cost)**



# Goal of TraceTracker

---

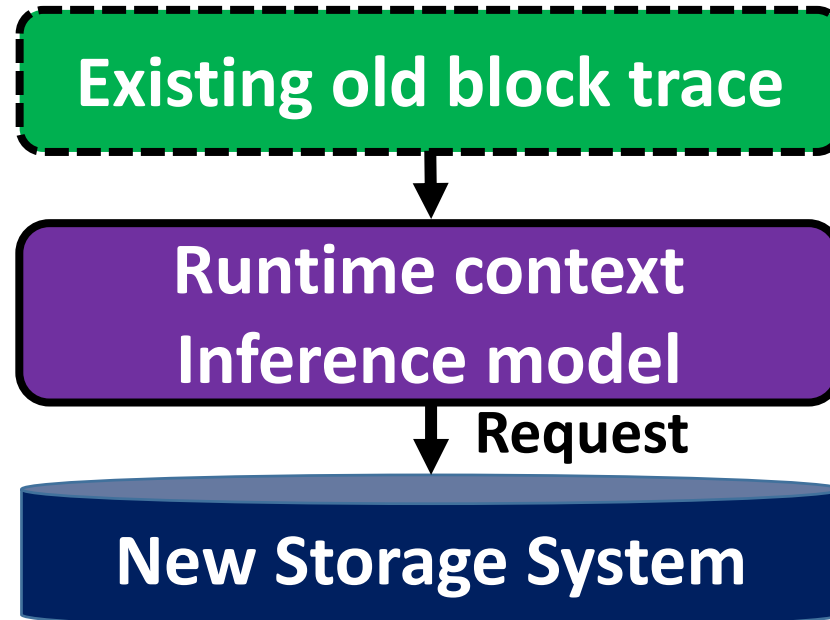


**Goal: Accurately reconstruct workloads without application execution (at low cost)**



# Goal of TraceTracker

---



**However, how runtime context can be inferred from block trace?**



# TraceTracker

---

1. Background of Block Trace

2. Trace Reconstruction Methods

**3. Insights on Trace Reconstruction**

4. TraceTracker Method

5. Evaluation





# Insights of TraceTracker method

---

**Goal:** Infer runtime contexts from timing information (Inter-arrival times) of old block traces



# Insights of TraceTracker method

---

**Goal:** Infer runtime contexts from timing information (Inter-arrival times) of old block traces

**Insight1:** Inter-arrival times decomposition



# Insights of TraceTracker method

---

**Goal:** Infer runtime contexts from timing information (Inter-arrival times) of old block traces

**Insight1:** Inter-arrival times decomposition

**Insight2:** Three inter-arrival time types



# Insights of TraceTracker method

---

**Goal:** Infer runtime contexts from timing information (Inter-arrival times) of old block traces

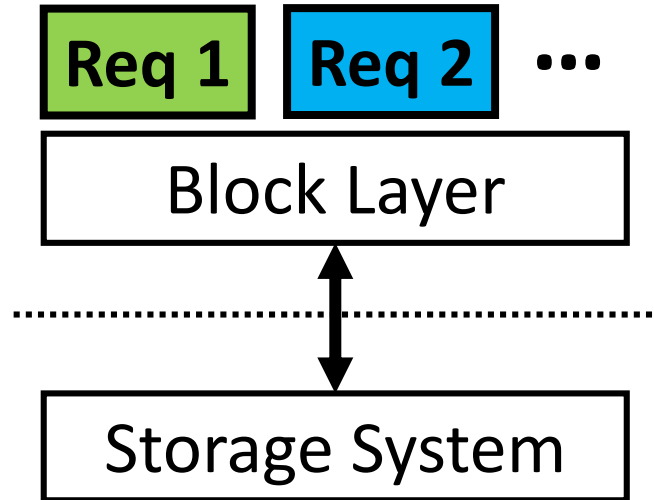
**Insight1:** Inter-arrival times decomposition

**Insight2:** Three inter-arrival time types

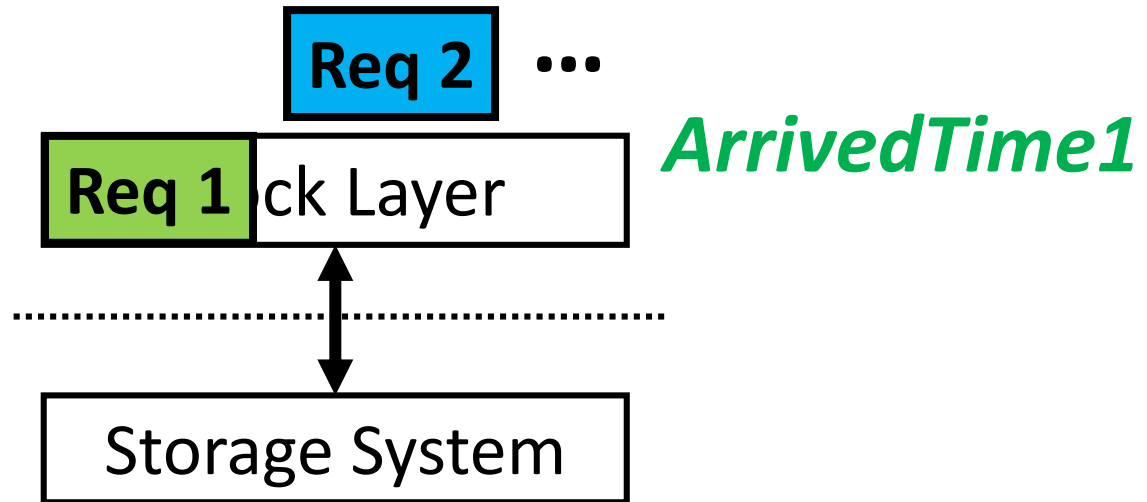
**Insight3:** Inter-arrival time components' characteristics



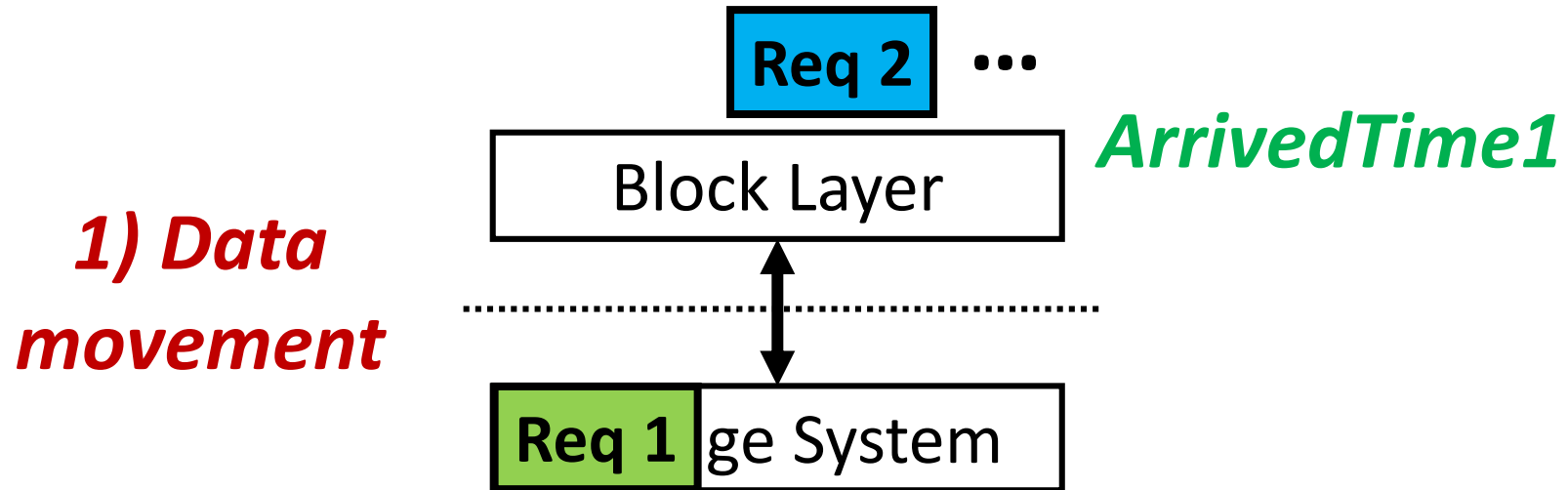
# Insight1: Inter-arrival times decomposition



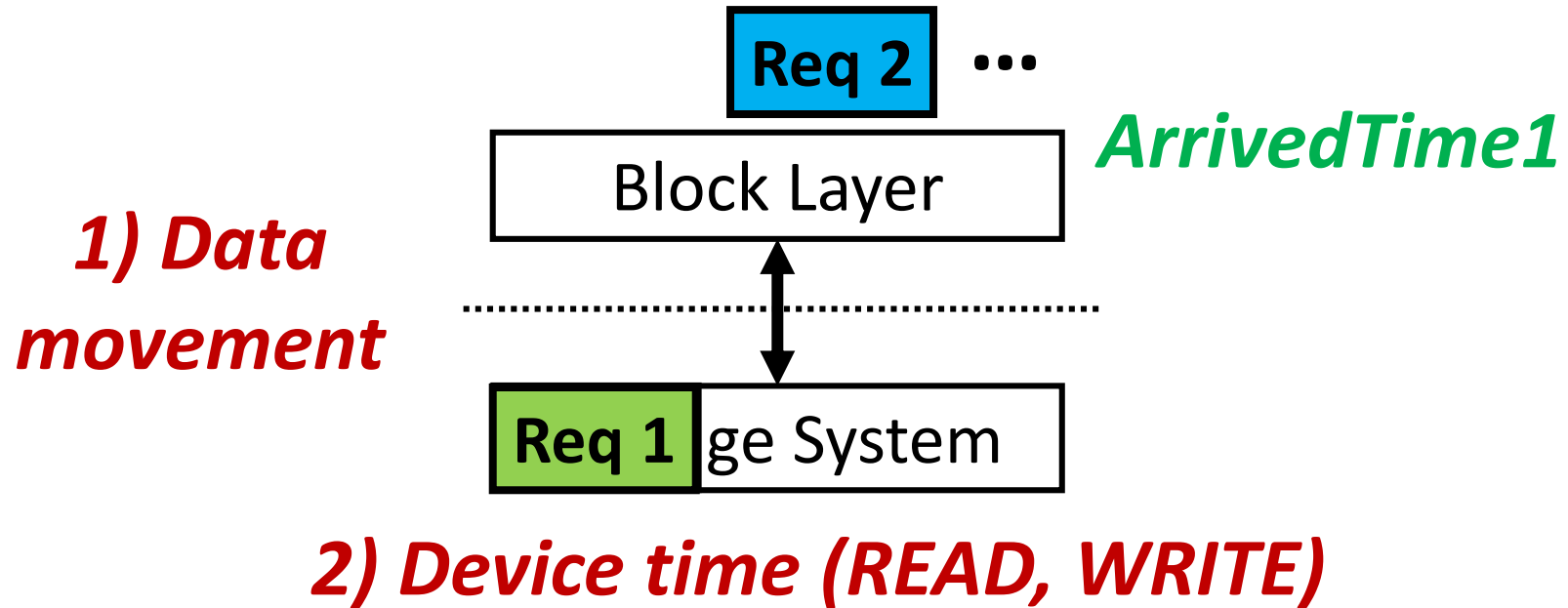
# Insight1: Inter-arrival times decomposition



# Insight1: Inter-arrival times decomposition

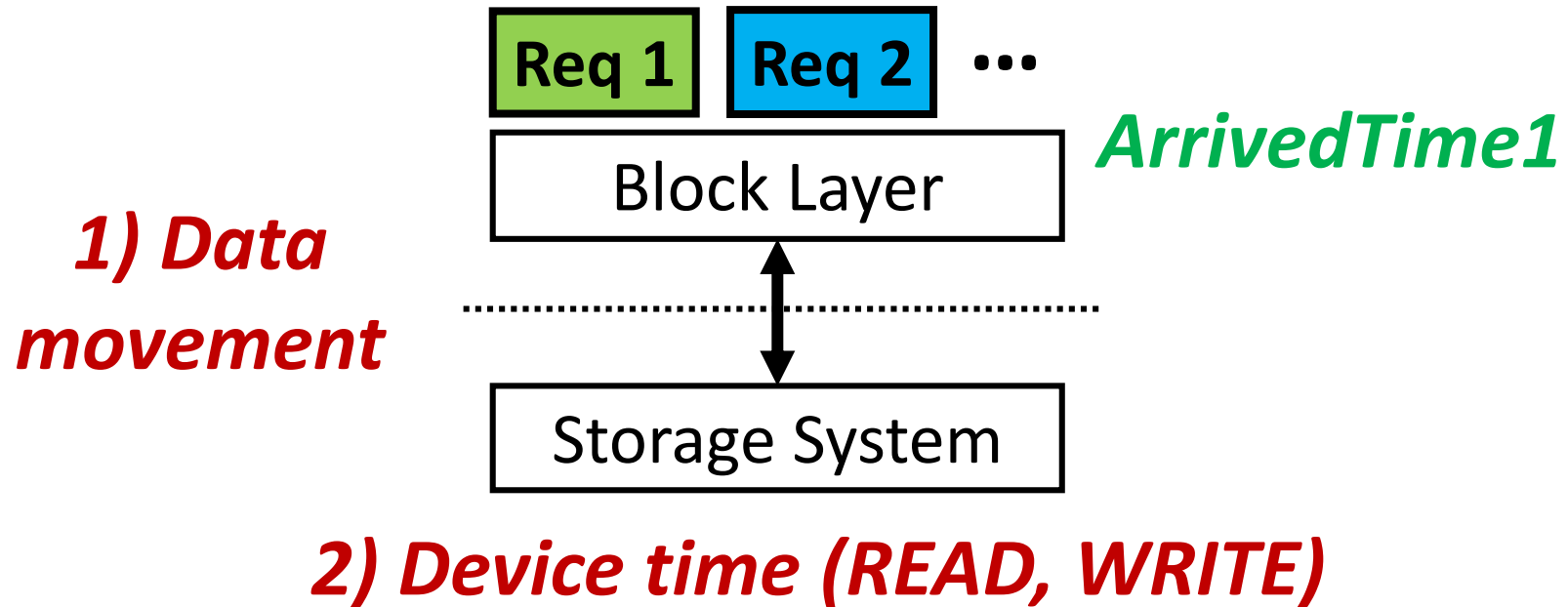


# Insight1: Inter-arrival times decomposition



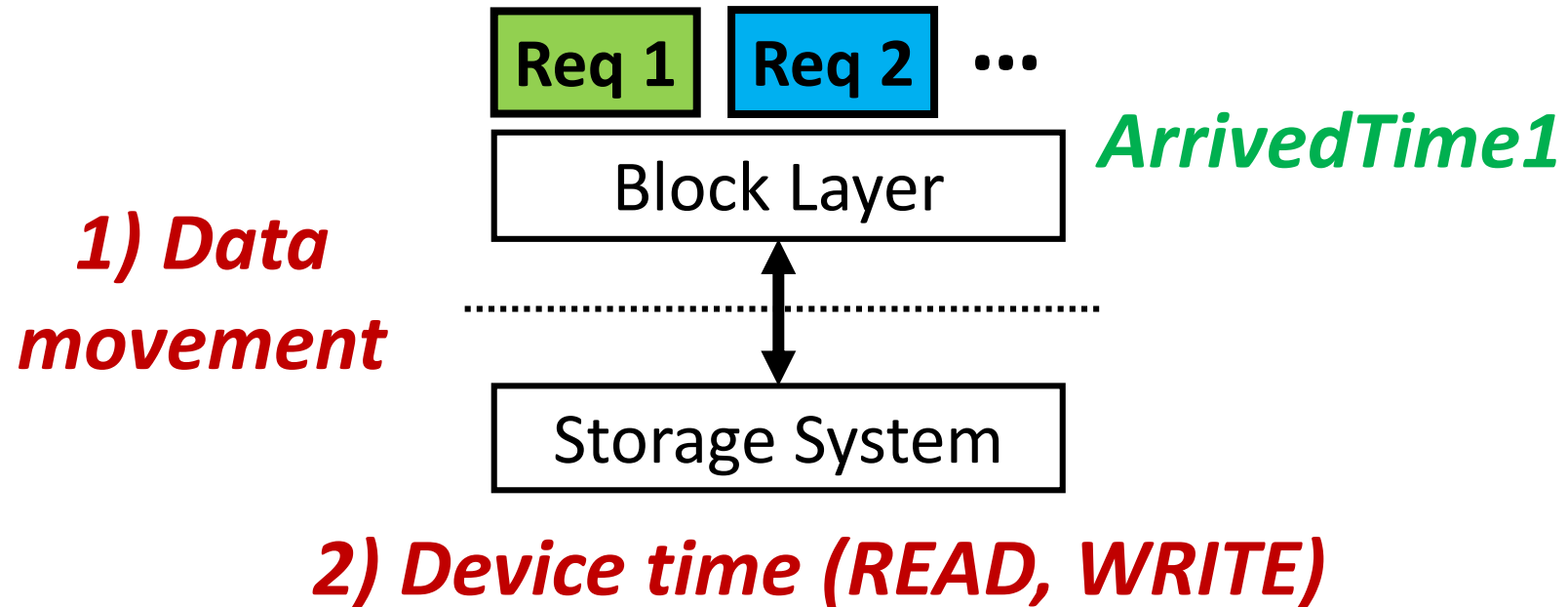


# Insight1: Inter-arrival times decomposition



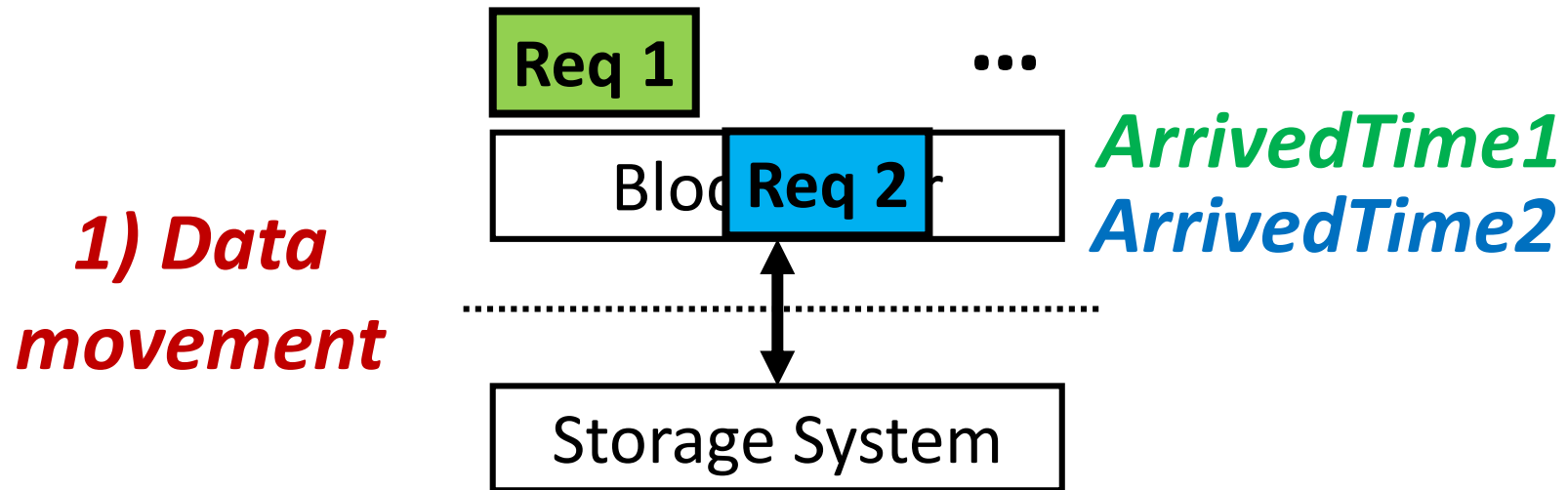
# Insight1: Inter-arrival times decomposition

## *3) User idle times*



# Insight1: Inter-arrival times decomposition

## *3) User idle times*

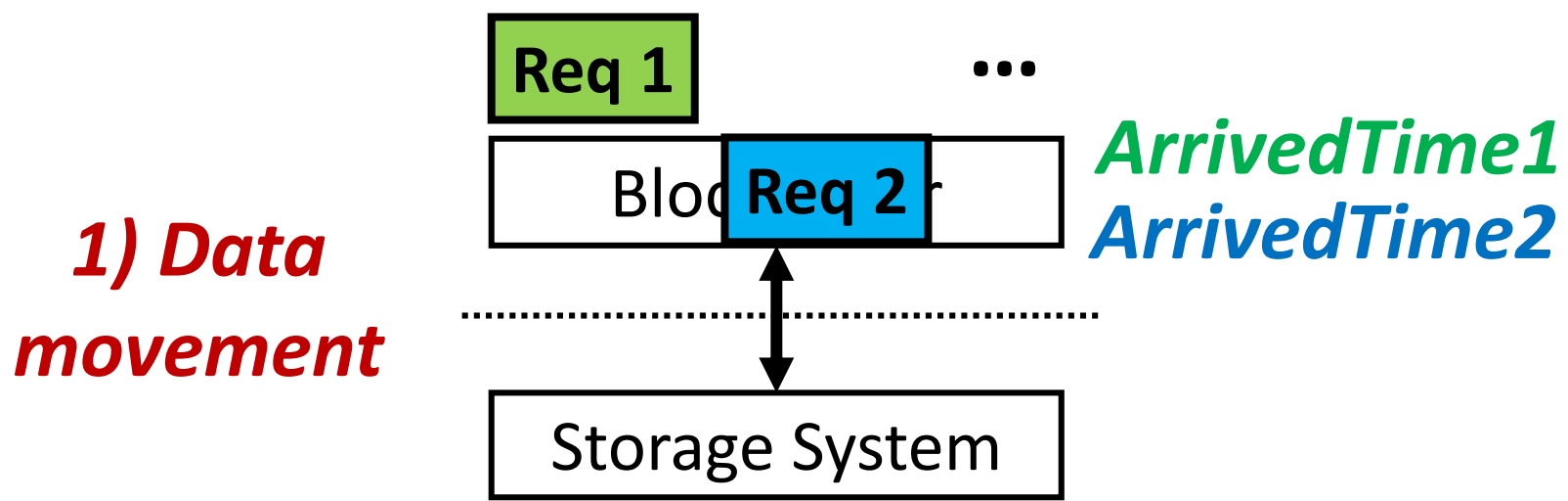


## *2) Device time (READ, WRITE)*



**Insight1:** Inter-arrival times decomposition

*3) User idle times*

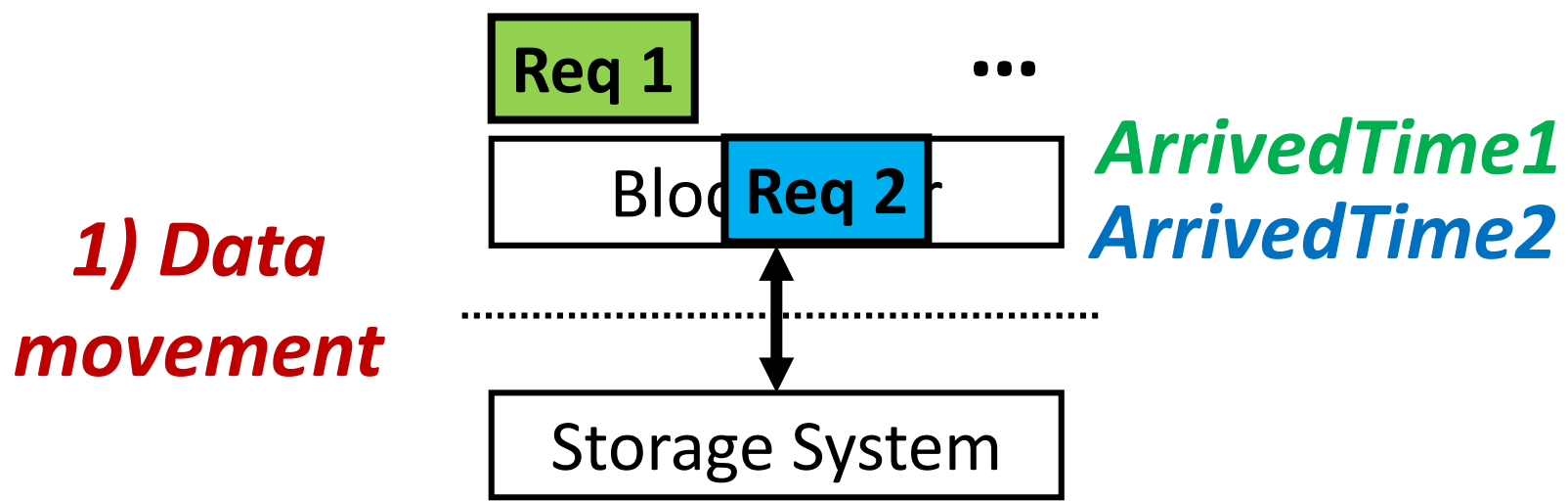


*2) Device time (READ, WRITE)*



**Insight1:** Inter-arrival times decomposition

*3) User idle times*

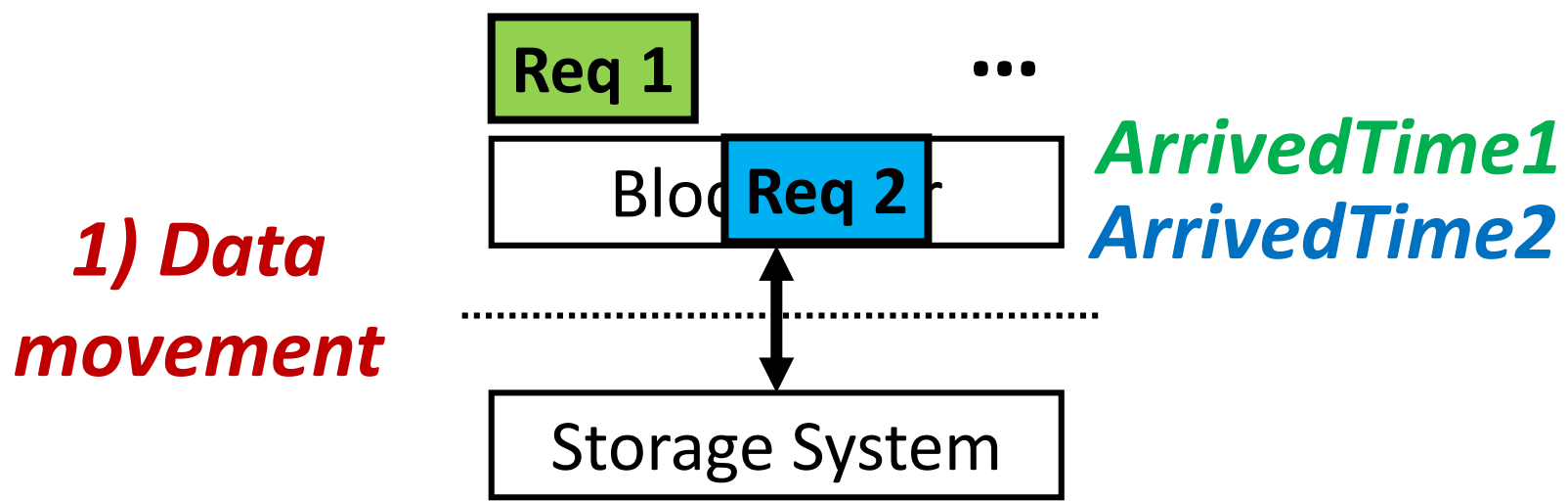


*2) Device time (READ, WRITE)*

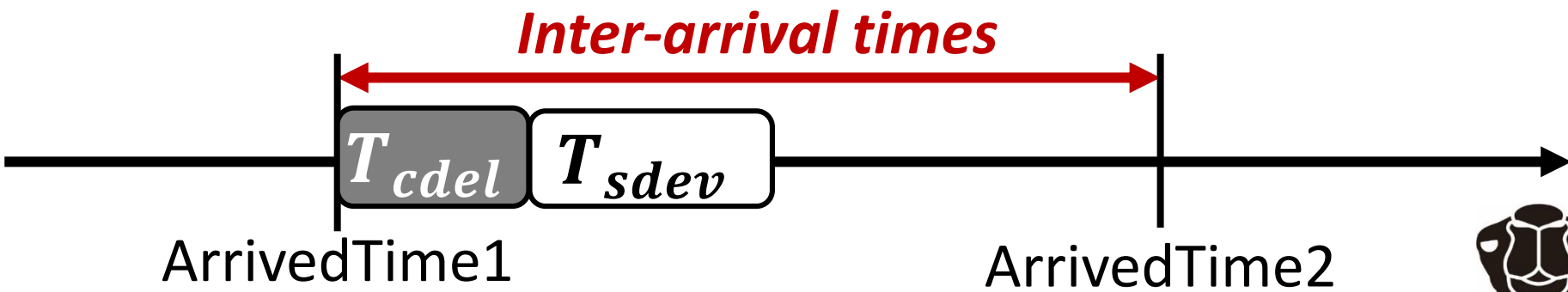


**Insight1:** Inter-arrival times decomposition

*3) User idle times*

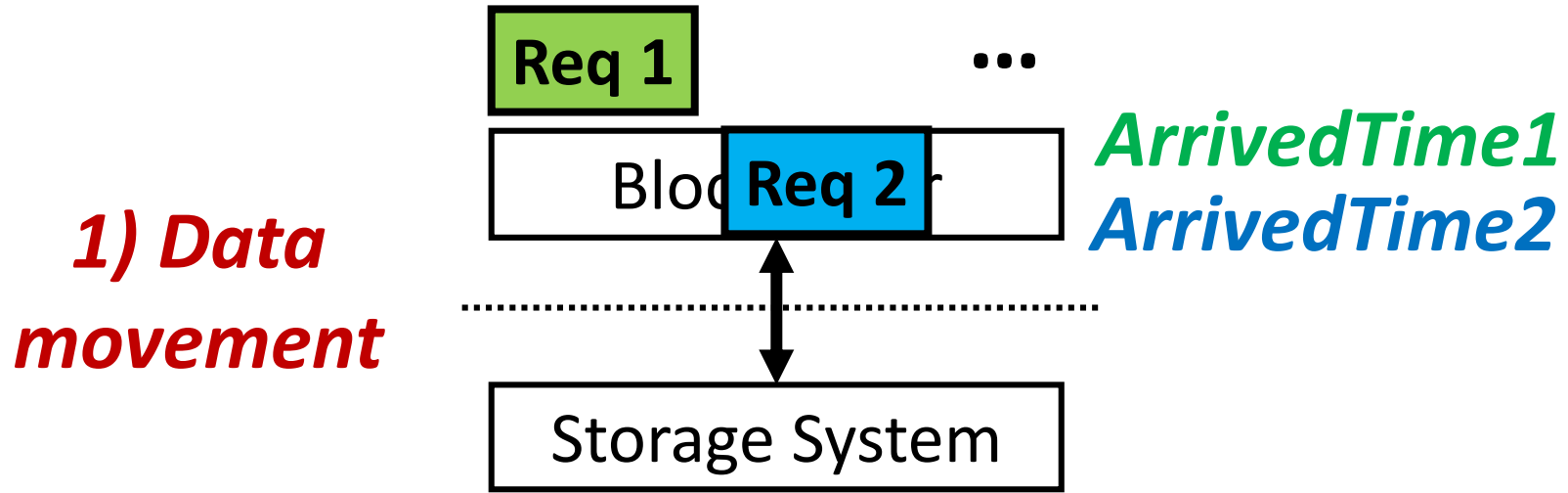


*2) Device time (READ, WRITE)*



**Insight1:** Inter-arrival times decomposition

*3) User idle times*



*2) Device time (READ, WRITE)*



**Insight2:** Three inter-arrival times types  
(I/O execution mode, application busy/idle status)





## Insight2: Three inter-arrival times types (I/O execution mode, application busy/idle status)

Sync I/O,  
Non-busy status



## Insight2: Three inter-arrival times types (I/O execution mode, application busy/idle status)

Sync I/O,  
Non-busy status



Sync I/O,  
Busy status



## Insight2: Three inter-arrival times types (I/O execution mode, application busy/idle status)

Sync I/O,  
Non-busy status




Sync I/O,  
Busy status



Async I/O



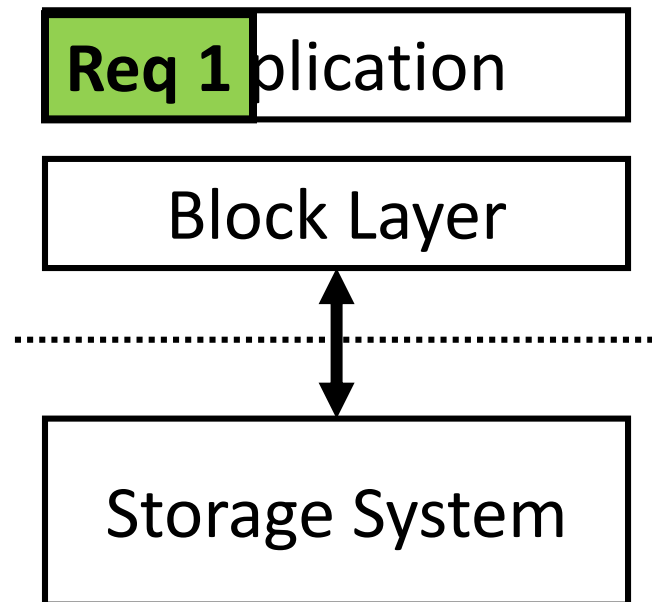
### Insight3: Different characteristics of three inter-arrival times components ( $T_{cdel}$ , $T_{sdev}$ , $T_{idle}$ )

	$T_{cdel}$	$T_{sdev}$	$T_{idle}$
Time length	 <i>short</i> <span style="float: right;"><i>long</i></span>		
Predictable	O	O	X (rare)



### Insight3: Different characteristics of three inter-arrival times components ( $T_{cdel}$ , $T_{sdev}$ , $T_{idle}$ )

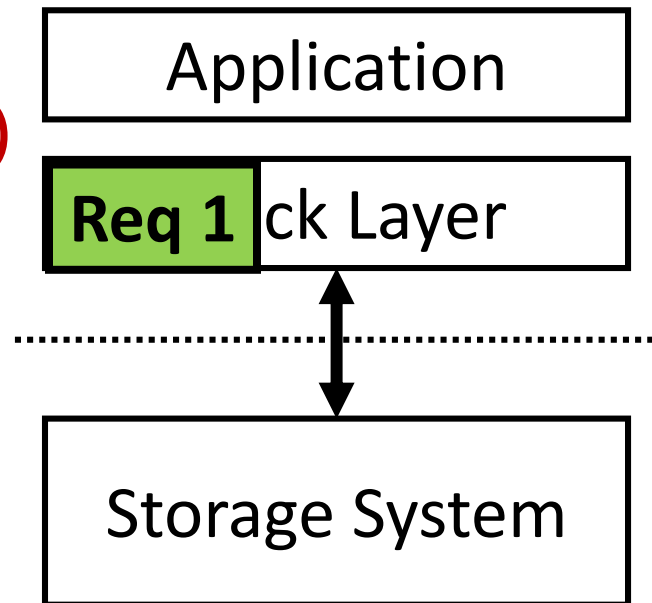
	$T_{cdel}$	$T_{sdev}$	$T_{idle}$
Time length	<div><div>←</div><div><math>\xrightarrow{\text{long}}</math></div></div> <div><i>short</i><span style="float: right;"><i>long</i></span></div>		
Predictable	O	O	X (rare)



### Insight3: Different characteristics of three inter-arrival times components ( $T_{cdel}$ , $T_{sdev}$ , $T_{idle}$ )

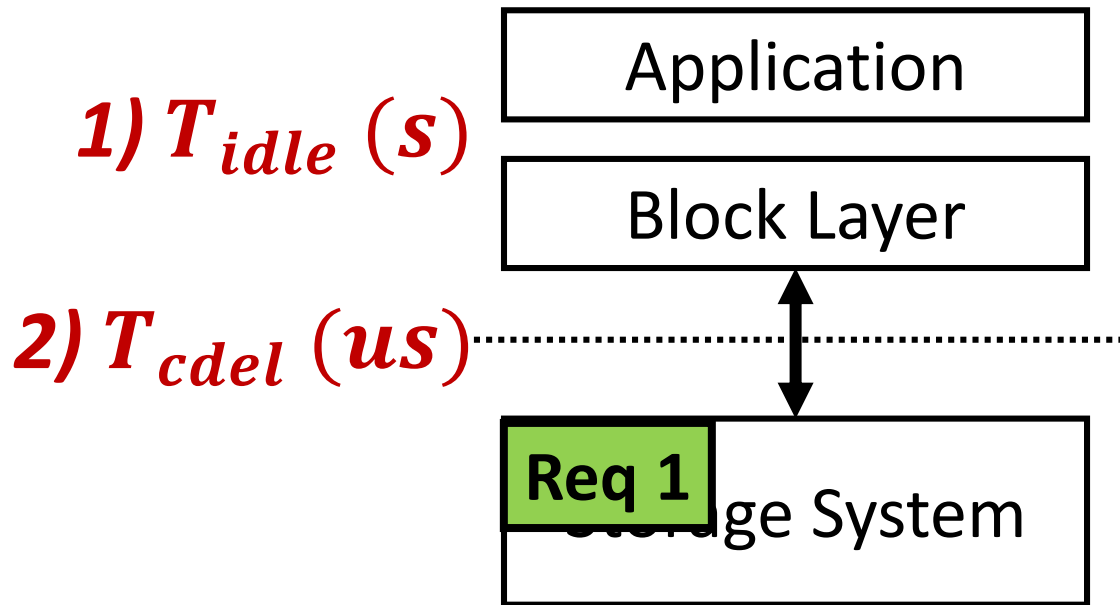
	$T_{cdel}$	$T_{sdev}$	$T_{idle}$
Time length	<div><div>←</div><div><math>\xrightarrow{\hspace{10em}}</math></div><div>→</div></div> <div><i>short</i><span style="float: right;"><i>long</i></span></div>		
Predictable	O	O	X (rare)

**1)  $T_{idle}$  (s)**



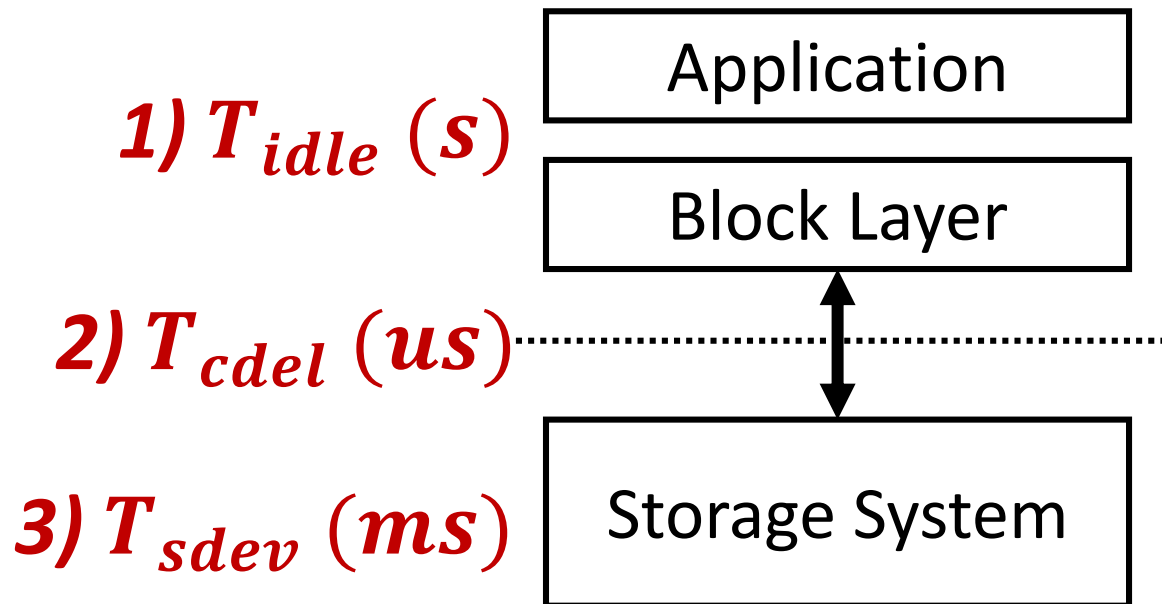
### Insight3: Different characteristics of three inter-arrival times components ( $T_{cdel}$ , $T_{sdev}$ , $T_{idle}$ )

	$T_{cdel}$	$T_{sdev}$	$T_{idle}$
Time length	<div><div>←</div><div></div><div>→</div></div> <div><i>short</i><span style="float: right;"><i>long</i></span></div>		
Predictable	O	O	X (rare)



# Insight3: Different characteristics of three inter-arrival times components ( $T_{cdel}$ , $T_{sdev}$ , $T_{idle}$ )

	$T_{cdel}$	$T_{sdev}$	$T_{idle}$
Time length	<div><div>←</div><div></div><div>→</div></div> <div><i>short</i><span style="float: right;"><i>long</i></span></div>		
Predictable	O	O	X (rare)





# TraceTracker

---

1. Background of Block Trace

2. Trace Reconstruction Methods

3. Insights on Trace Reconstruction

4. TraceTracker Method

5. Evaluation



# Overview of TraceTracker

---



# Overview of TraceTracker

---

*Pre-software  
evaluation*



# Overview of TraceTracker

---

*Pre-software  
evaluation*  *Hardware  
evaluation*



# Overview of TraceTracker

---



# Overview of TraceTracker

---

*Pre-software  
evaluation*

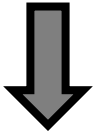


*Hardware  
evaluation*



*Post-software  
evaluation*

Old block traces

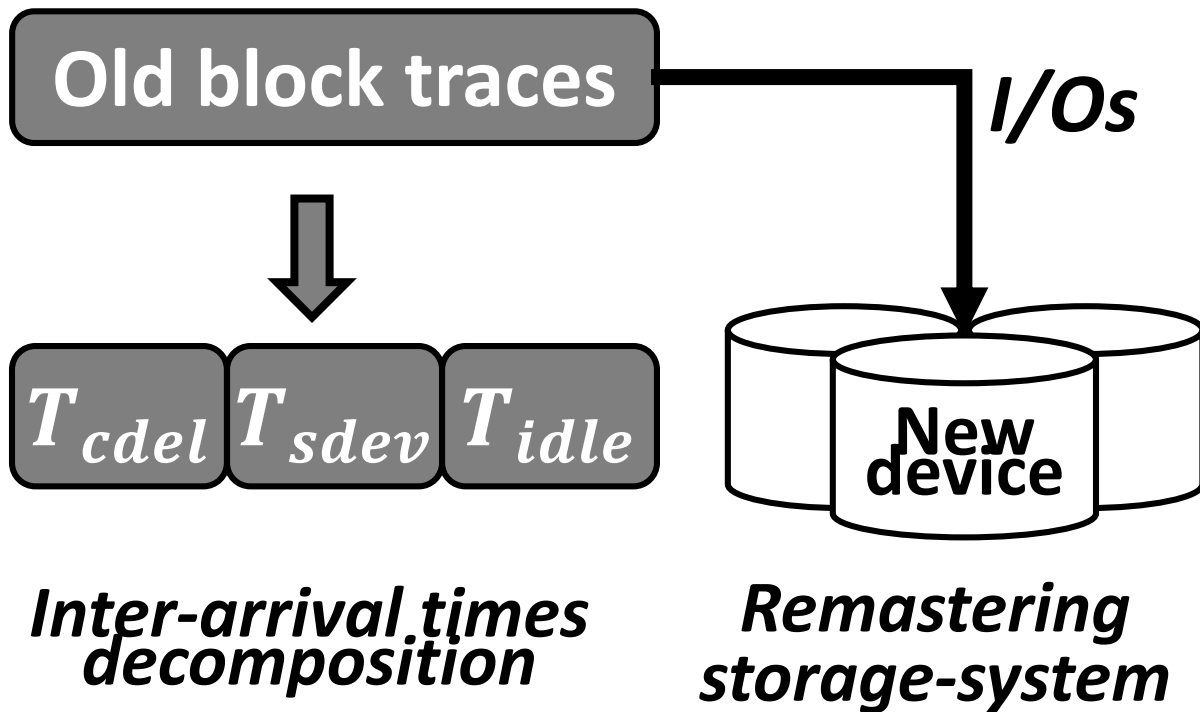


*Inter-arrival times  
decomposition*



# Overview of TraceTracker

*Pre-software evaluation* → *Hardware evaluation* → *Post-software evaluation*



# Overview of TraceTracker

*Pre-software  
evaluation*

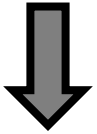


*Hardware  
evaluation*



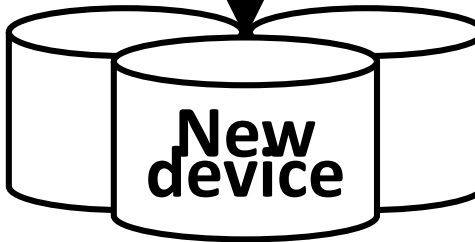
*Post-software  
evaluation*

Old block traces



*Inter-arrival times  
decomposition*

$I/Os$



*Remastering  
storage-system*





# Overview of TraceTracker

*Pre-software  
evaluation*

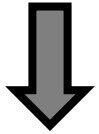


*Hardware  
evaluation*



*Post-software  
evaluation*

**Old block traces**

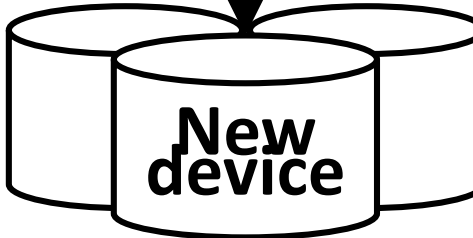


$T_{c\text{del}}$   $T_{s\text{dev}}$   $T_{i\text{idle}}$

*Inter-arrival times  
decomposition*



*I/Os*



**New  
device**

*Remastering  
storage-system*

**Replayed traces**



**New block traces**

*I/O execution mode  
Inference model*



# Overview of TraceTracker

*Pre-software  
evaluation*

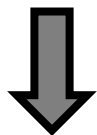


*Hardware  
evaluation*



*Post-software  
evaluation*

**Old block traces**

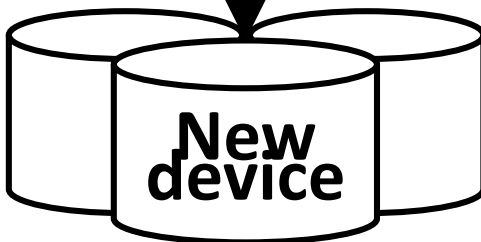


$T_{cdel}$   $T_{sdev}$   $T_{idle}$

*Inter-arrival times  
decomposition*



*I/Os*



**Remastering  
storage-system**

**Replayed traces**

*Mode flag*



**New block traces**

*I/O execution mode  
Inference model*



# Overview of TraceTracker

*Pre-software  
evaluation*

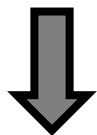


*Hardware  
evaluation*



*Post-software  
evaluation*

**Old block traces**

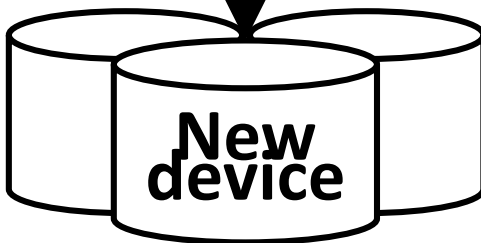


$T_{cdel}$   $T_{sdev}$   $T_{idle}$

*Inter-arrival times  
decomposition*



*I/Os*



*Remastering  
storage-system*

**Replayed traces**

*Mode flag*

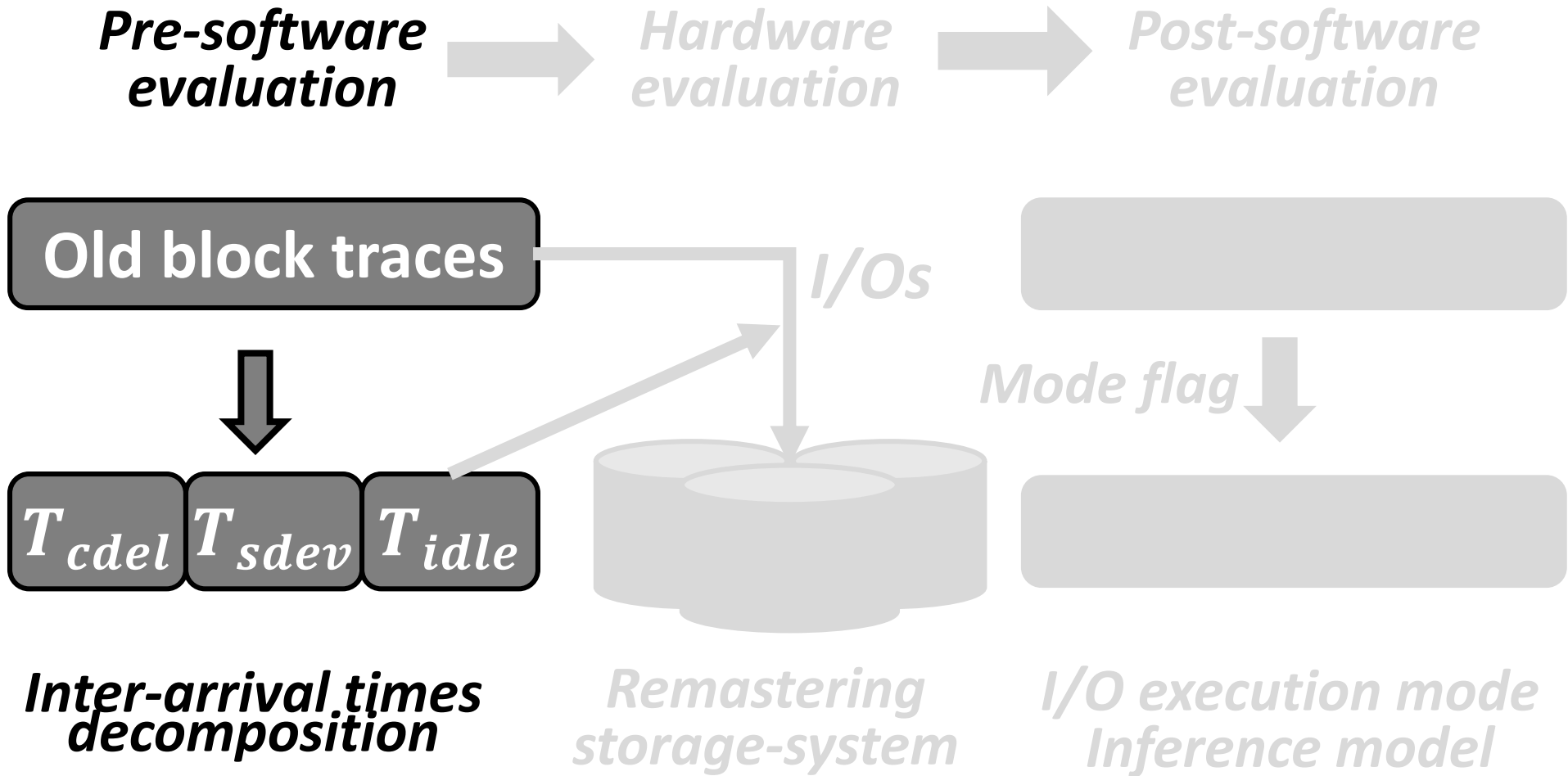


**New block traces**

*I/O execution mode  
Inference model*



# Overview of TraceTracker



# Request classification

---

**Key idea:** Classify all requests in trace with access pattern, operation type, request length

Old block traces

Req1: Read Len1 ..  
Req2: Read Len2 ..  
Req3: Write Len3 ..  
...  
ReqN: Write LenN ..



# Request classification

---

**Key idea:** Classify all requests in trace with access pattern, operation type, request length

Old block traces

Req1: Read Len1 ..  
Req2: Read Len2 ..  
Req3: Write Len3 ..  
...  
ReqN: Write LenN ..



Length1

SeqRead ..

SeqWrite ..

RandRead ..

RandWrite ..

...

Length L

SeqRead ..

SeqWrite ..

RandRead ..

RandWrite ..



# Request classification

---

**Key idea:** Classify all requests in trace with access pattern, operation type, request length

Old block traces

Req1: Read Len1 ..  
Req2: Read Len2 ..  
Req3: Write Len3 ..  
...  
ReqN: Write LenN ..



Length1

SeqRead ..

SeqWrite ..

RandRead ..

RandWrite ..

...

Length L

SeqRead ..

SeqWrite ..

RandRead ..

RandWrite ..



# Inter-arrival times distribution

---

**Key idea:** Easily decompose  $T_{intt}$  from same category requests'  $T_{intt}$  distribution





# Inter-arrival times distribution

---

**Key idea:** Easily decompose  $T_{\text{intt}}$  from same category requests'  $T_{\text{intt}}$  distribution

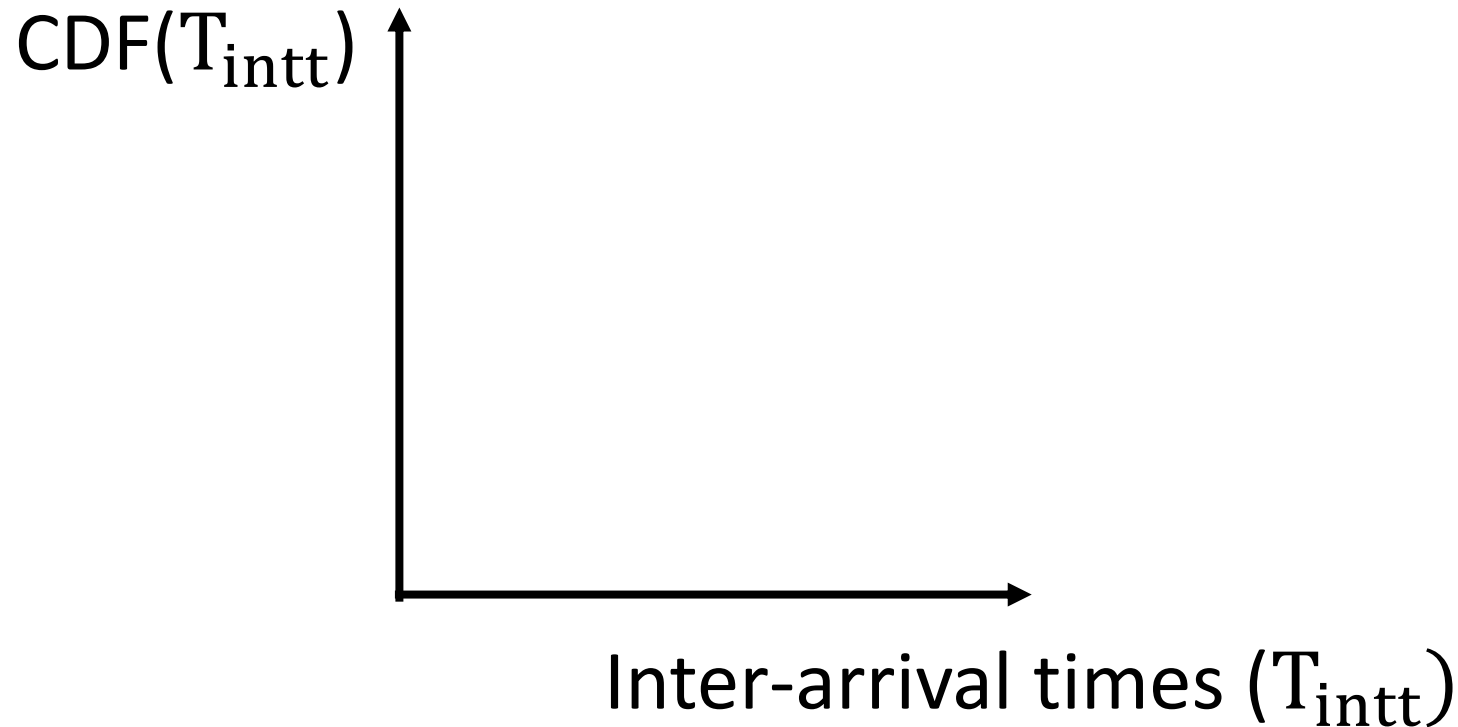
Length1 sequential read



# Inter-arrival times distribution

---

**Key idea:** Easily decompose  $T_{\text{intt}}$  from same category requests'  $T_{\text{intt}}$  distribution



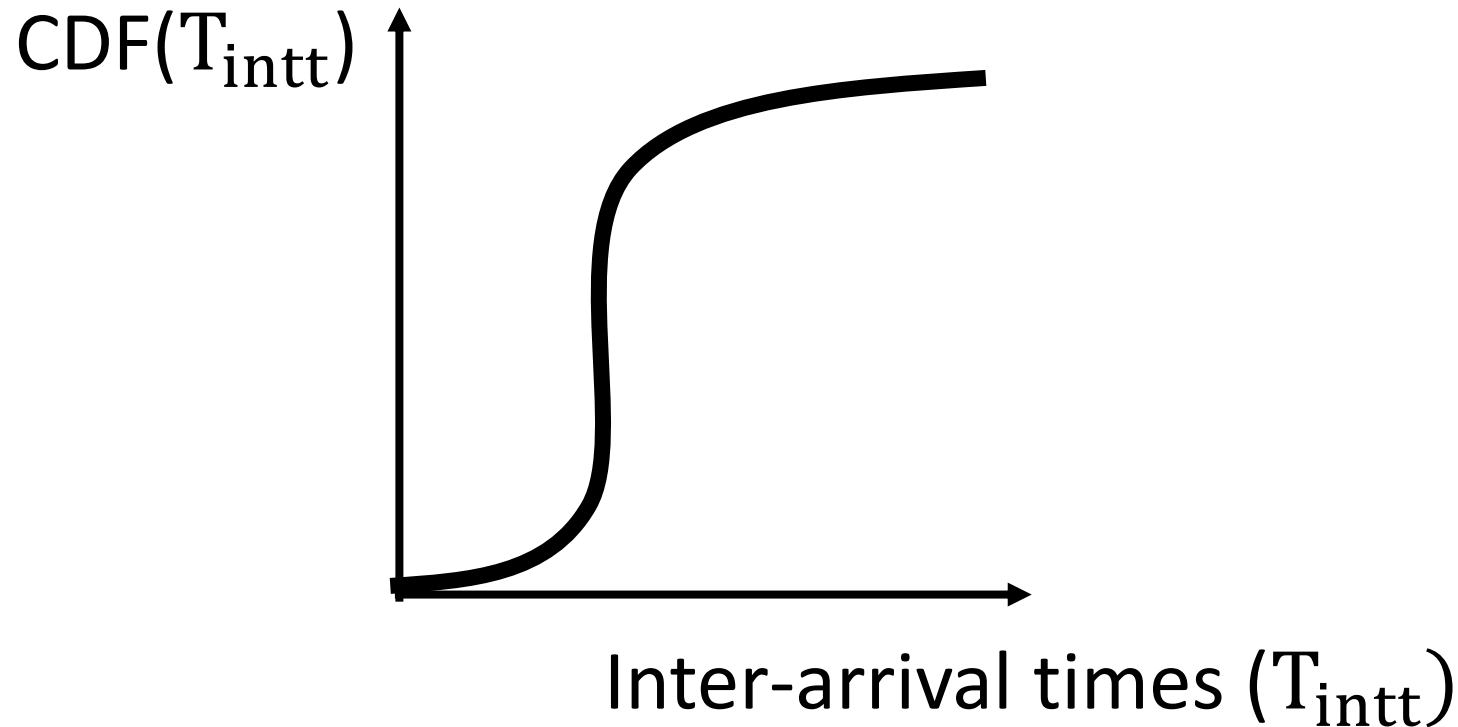
Length1 sequential read



# Inter-arrival times distribution

---

**Key idea:** Easily decompose  $T_{\text{intt}}$  from same category requests'  $T_{\text{intt}}$  distribution



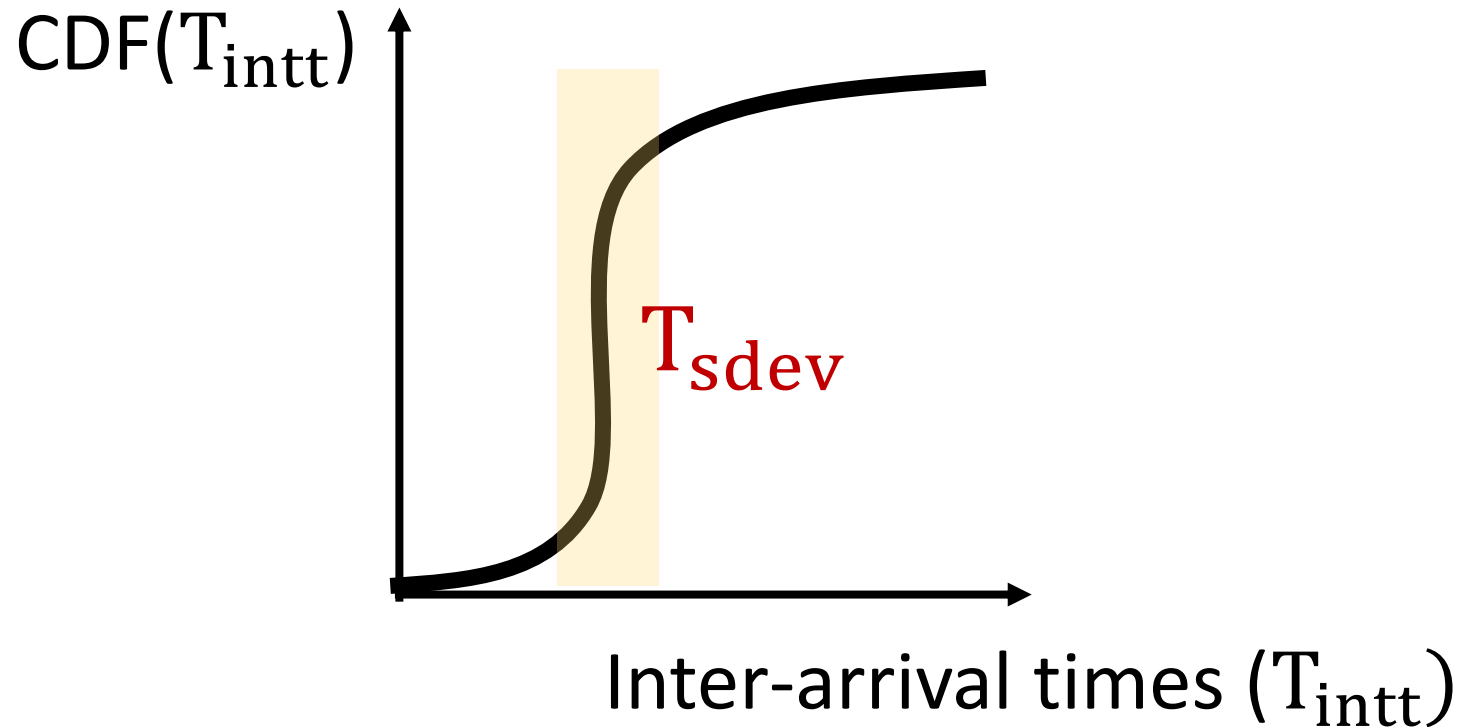
Length1 sequential read



# Inter-arrival times distribution

---

**Key idea:** Easily decompose  $T_{\text{intt}}$  from same category requests'  $T_{\text{intt}}$  distribution



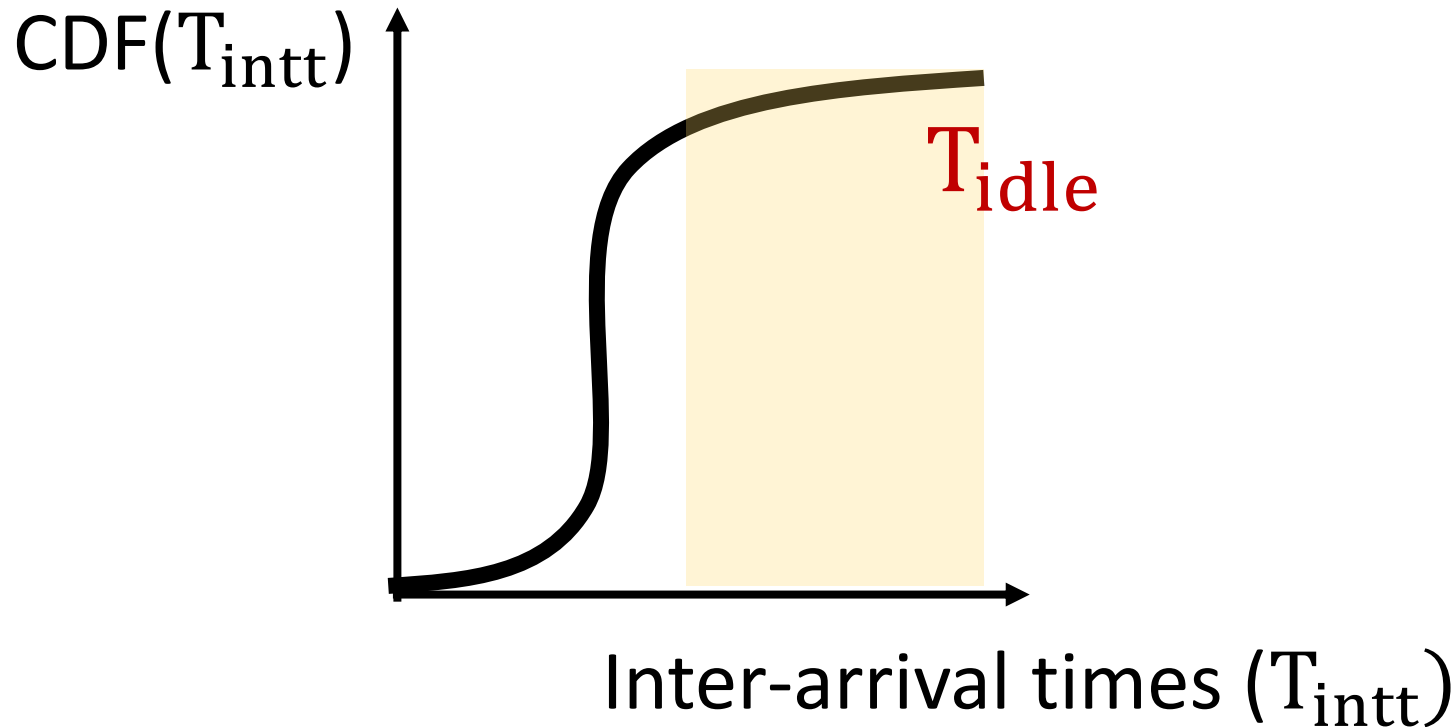
Length1 sequential read



# Inter-arrival times distribution

---

**Key idea:** Easily decompose  $T_{\text{intt}}$  from same category requests'  $T_{\text{intt}}$  distribution



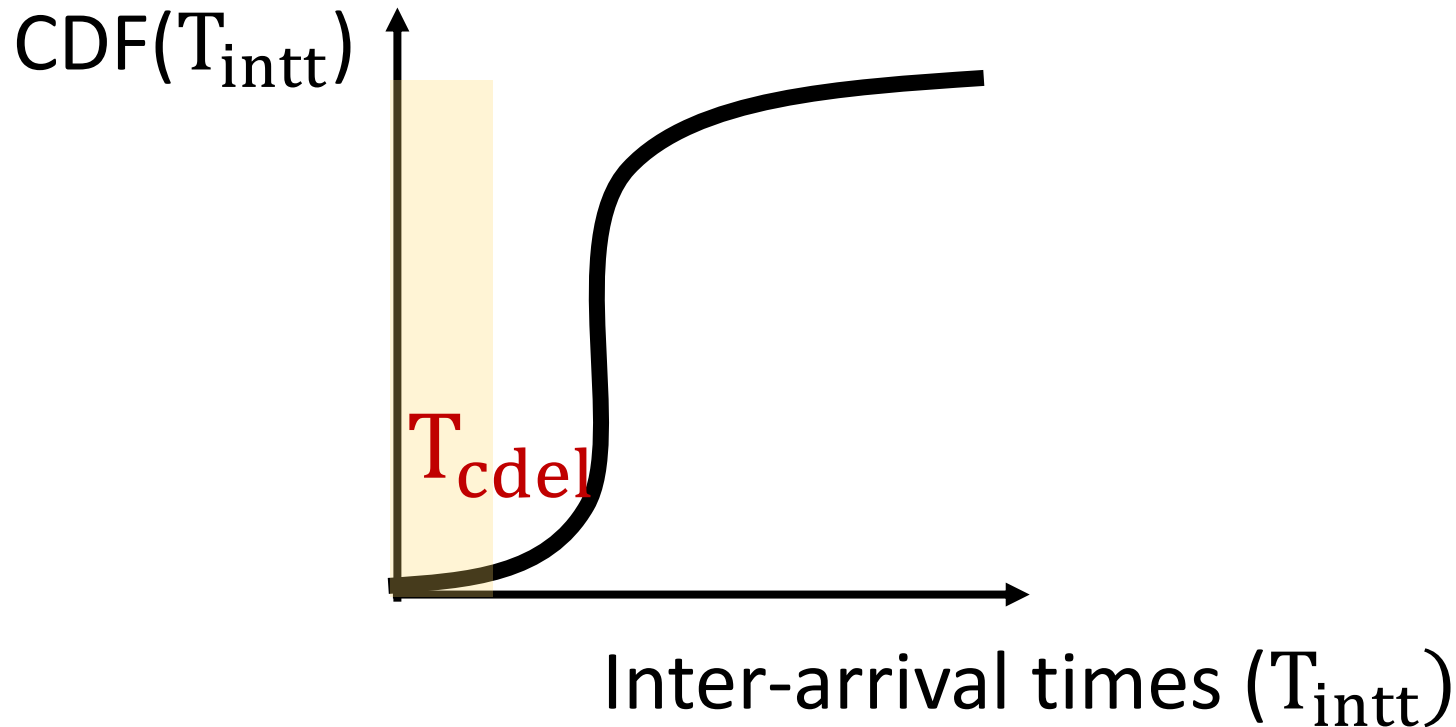
Length1 sequential read



# Inter-arrival times distribution

---

**Key idea:** Easily decompose  $T_{\text{intt}}$  from same category requests'  $T_{\text{intt}}$  distribution



Length1 sequential read



# Overview of TraceTracker

*Pre-software  
evaluation*



*Hardware  
evaluation*



*Post-software  
evaluation*

**Old block traces**



$T_{idle}$

*Inter-arrival times  
decomposition*



*I/Os*



*Remastering  
storage-system*

*Mode flag*



*I/O execution mode  
Inference model*



# Trace replay

---

**In-house application**





# Trace replay

---

## Old block trace

Req1: READ 1513413400 1024  
Req2: WRITE 1513413800 2048  
...



*T<sub>idle</sub>*



**In-house application**



# Trace replay

---

## Old block trace

Req1: READ 1513413400 1024  
Req2: WRITE 1513413800 2048  
...

$T_{idle}$

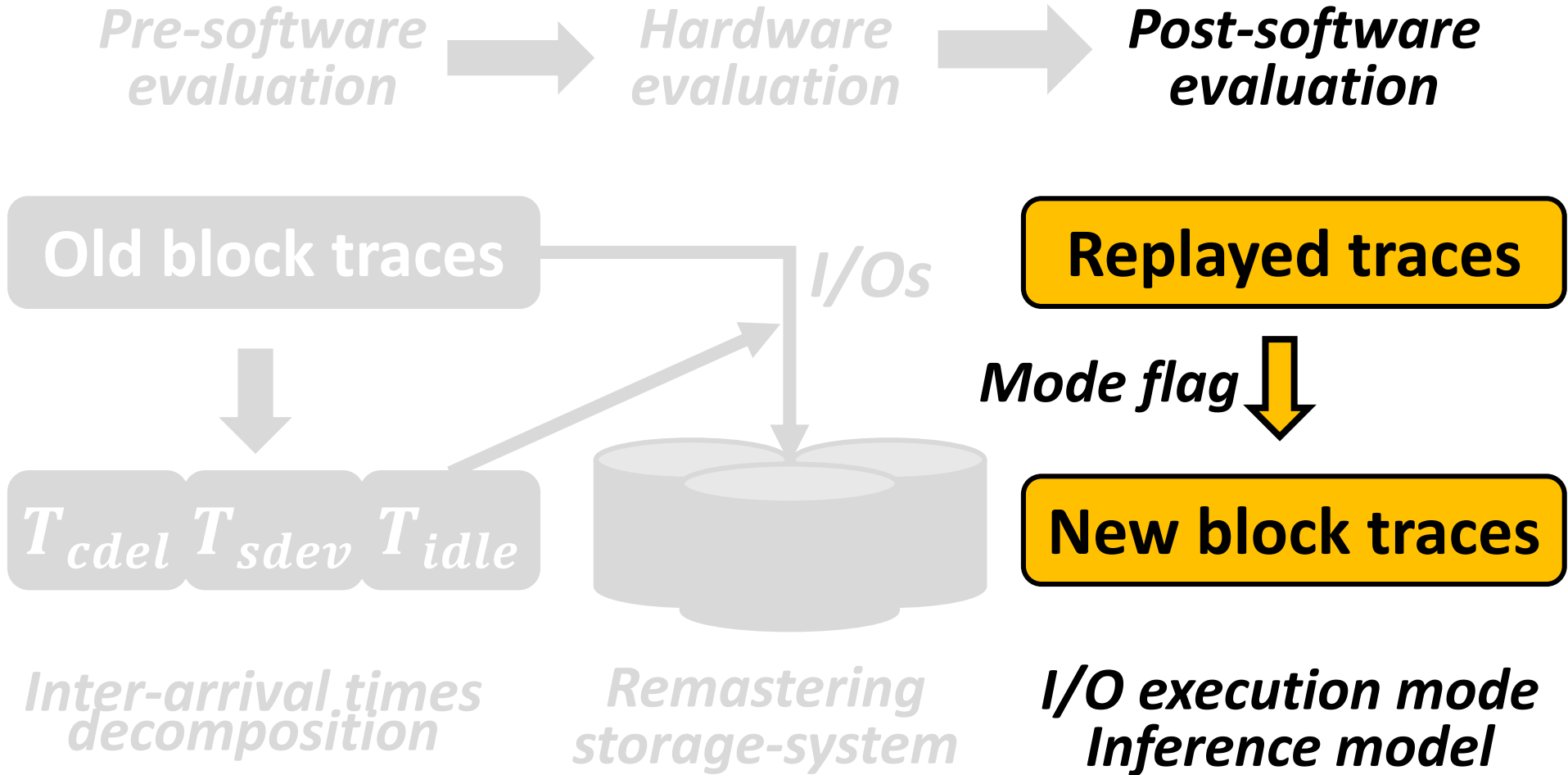
In-house application

I/O

New device



# Overview of TraceTracker



# I/O execution mode inference

**Key idea:** Asynchronous I/O has short period which is less than storage I/O latency ( $T_{sdev}$ )

Old traces

$T_{intt}$	1
$T_{intt}$	2
$T_{intt}$	3
...	

>  
>  
<

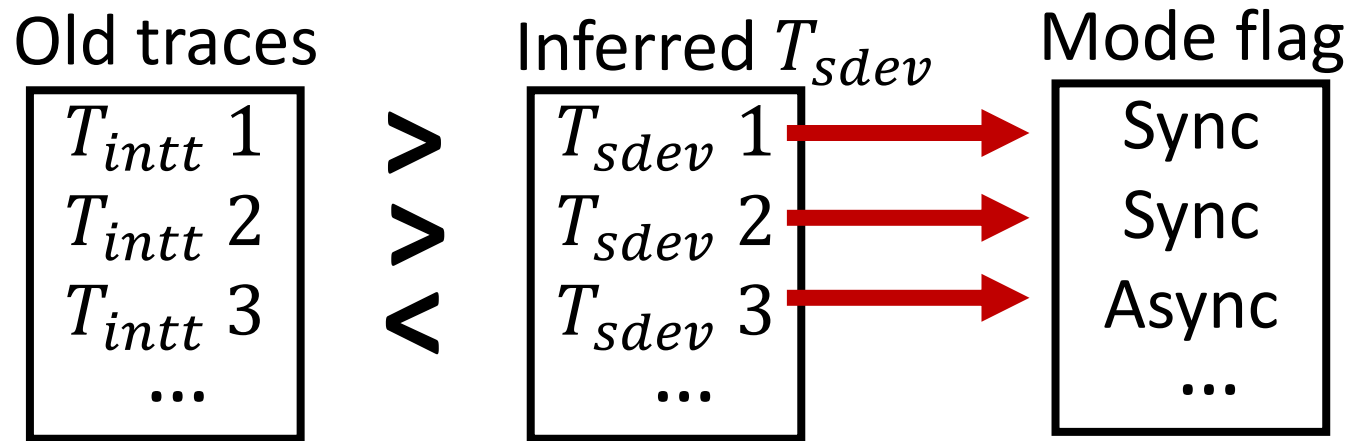
Inferred  $T_{sdev}$

$T_{sdev}$	1
$T_{sdev}$	2
$T_{sdev}$	3
...	



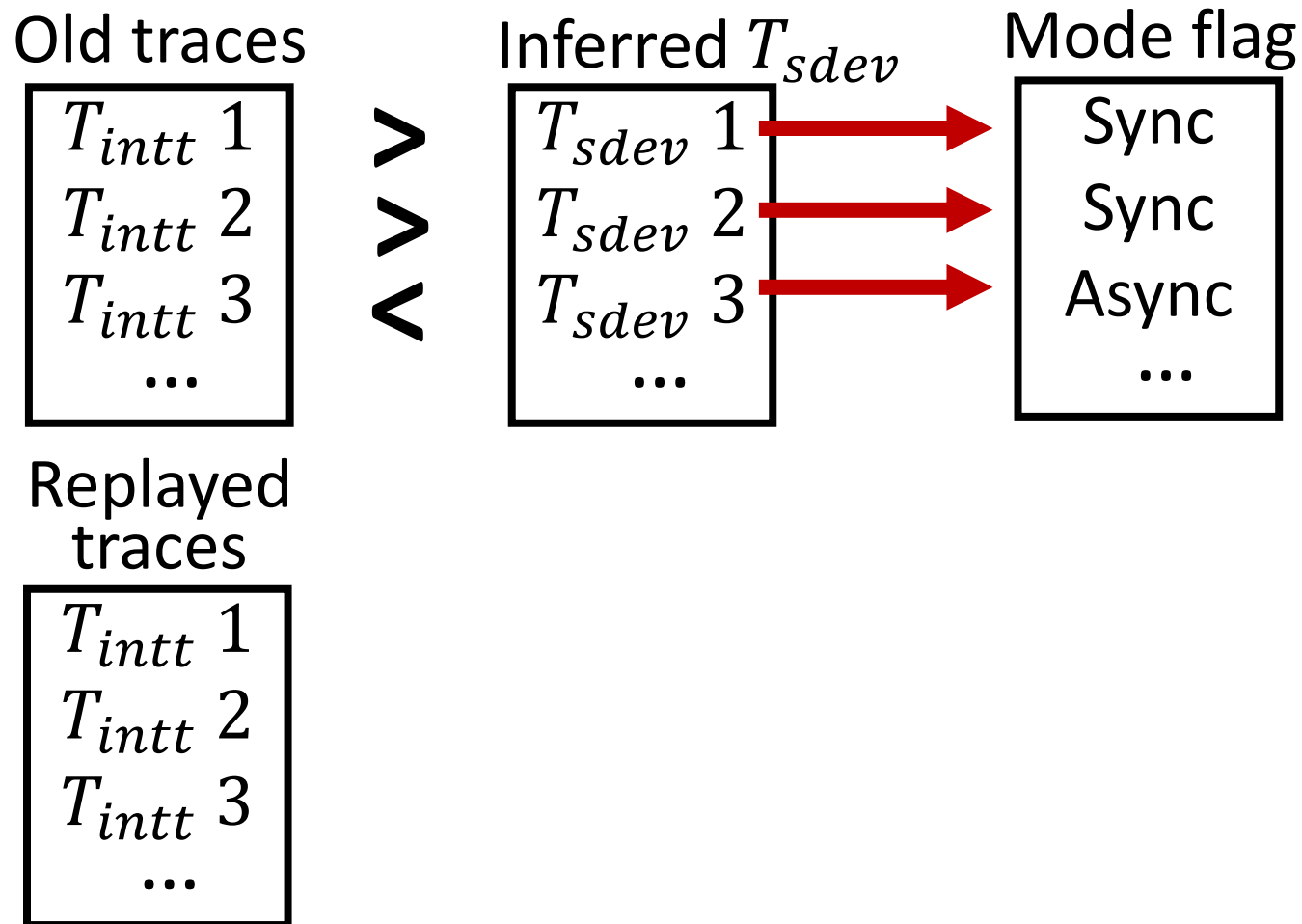
# I/O execution mode inference

**Key idea:** Asynchronous I/O has short period which is less than storage I/O latency ( $T_{sdev}$ )



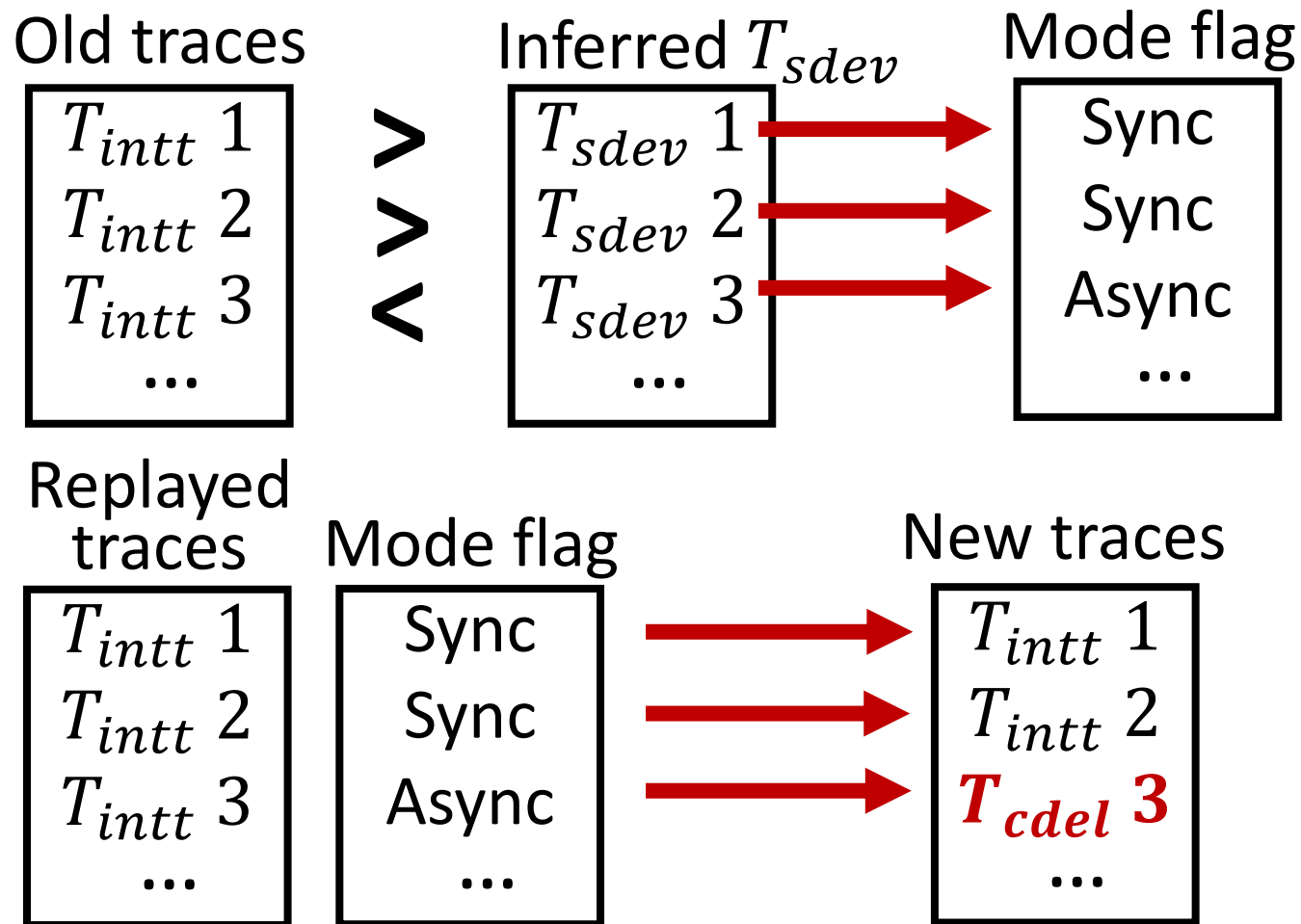
# I/O execution mode inference

**Key idea:** Asynchronous I/O has short period which is less than storage I/O latency ( $T_{sdev}$ )



# I/O execution mode inference

**Key idea:** Asynchronous I/O has short period which is less than storage I/O latency ( $T_{sdev}$ )



# TraceTracker

---

1. Background of Block Trace

2. Trace Reconstruction Methods

3. Insights on Trace Reconstruction

4. TraceTracker Method

5. Evaluation





# Evaluation Methodology

---

- **Evaluation node**

- All-flash array (by grouping four NVM Express SSDs)

- **Target block traces**

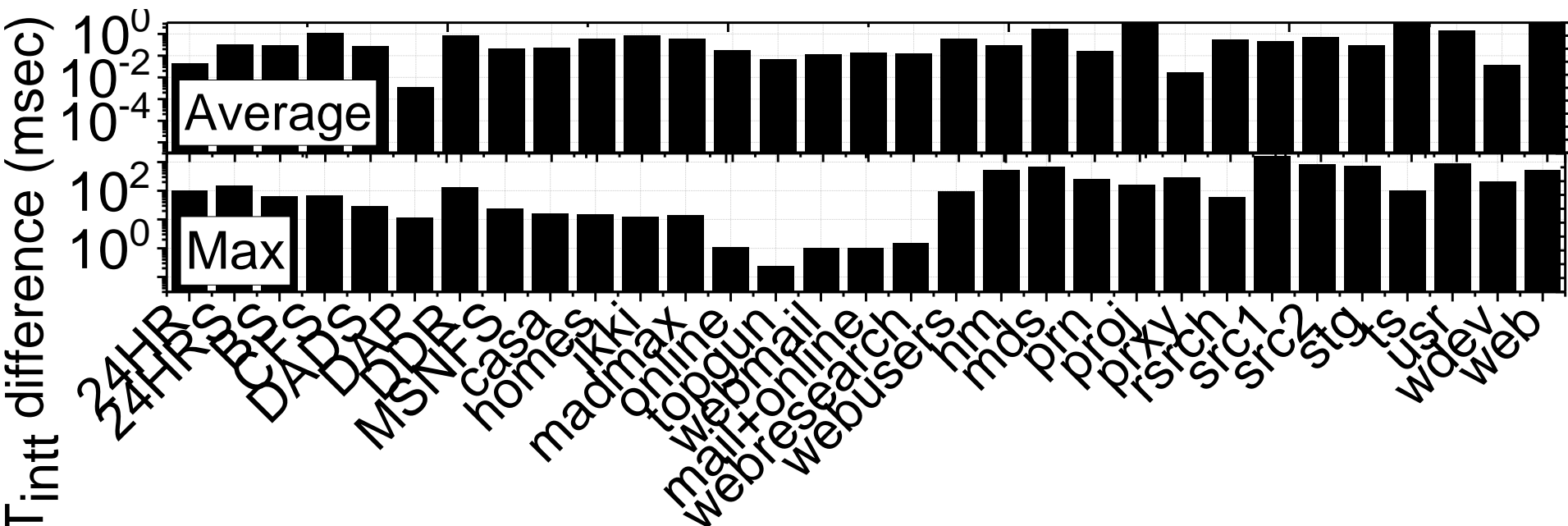
- 577 of large-scale block traces (around 2007~2009)

- **Reconstruction techniques**

- **Acceleration**: Reconstruction by shortening  $T_{intt}$
- **Revision**: Replaying block trace on all-flash array
- **Fixed-th**: An advanced revision method by inferring  $T_{idle}$  with **a fixed threshold**
- **Dynamic**: Reconstructions using our inference model, but **with no post-processing**.
- **TraceTracker**: Hardware/software co-evaluation



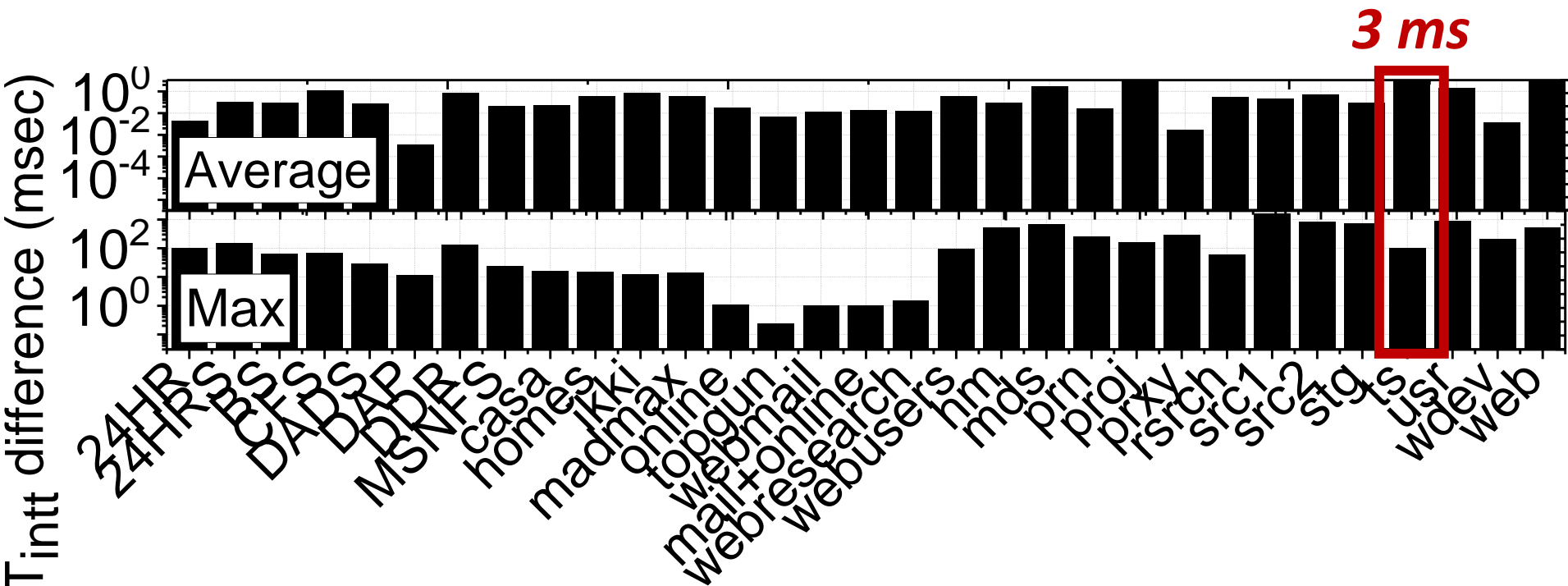
# Evaluation #1: Different I/O timing



- New traces have, on average, 0.677 ms **shorter** inter-arrival time than old block traces
- Ex) 'ts' has 3 ms shorter  $T_{intt}$ , on average.



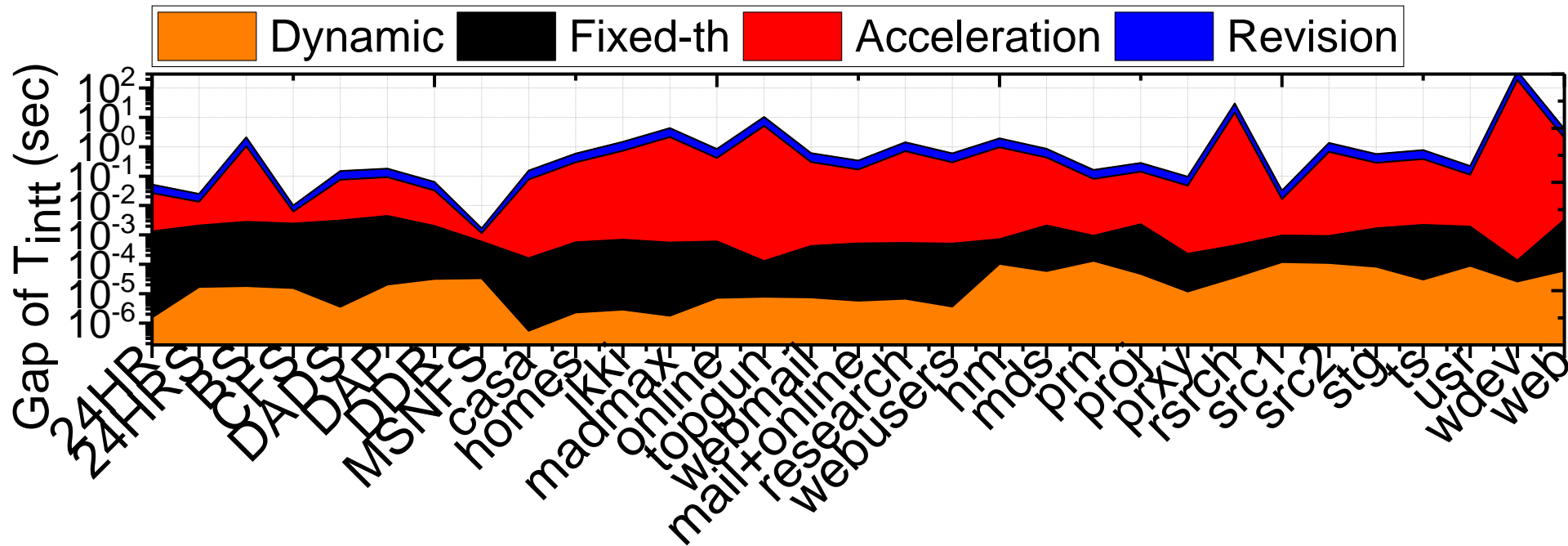
# Evaluation #1: Different I/O timing



- New traces have, on average, 0.677 ms **shorter** inter-arrival time than old block traces
- Ex) 'ts' has 3 ms shorter  $T_{intt}$ , on average.



# Evaluation #2: Inaccuracy of reconstruction

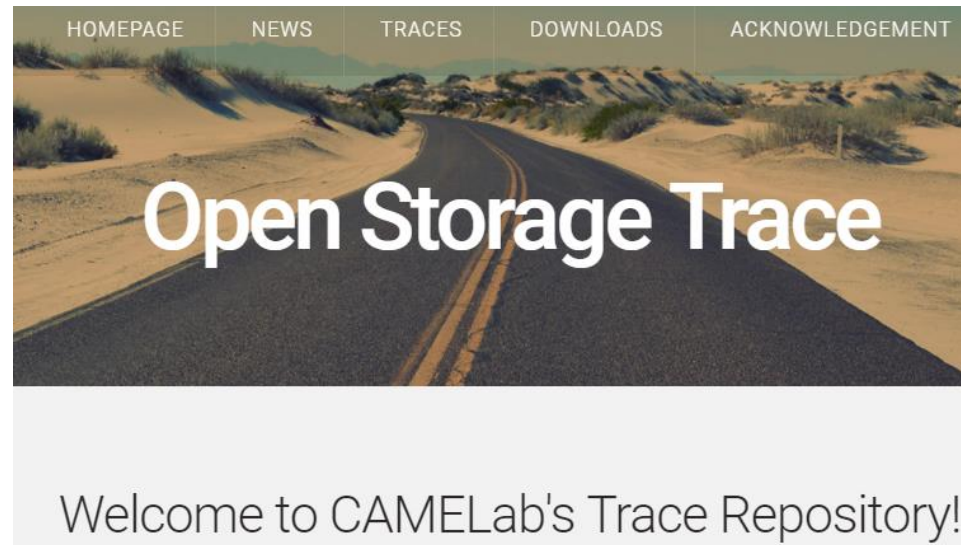


- ***Acceleration*** and ***Revision***: 7.08, 7.15 sec due to no information for  $T_{idle}$
- ***Fixed-th*** : 1.3 ms (inaccurate  $T_{idle}$ )
- ***Dynamic*** : 0.035 ms (inaccurate IO-mode)



# Conclusion

- ***TraceTracker*** is trace reconstruction method which remasters the storage latency while maintaining runtime contexts of target traces.
  - **We're preparing brand new traces for open-license!**
- Trace download is available at ***[trace.camelab.org](https://trace.camelab.org)***



**KANDEMIR : All-Flash Array based HPC Testbed**



# *TraceTracker:*

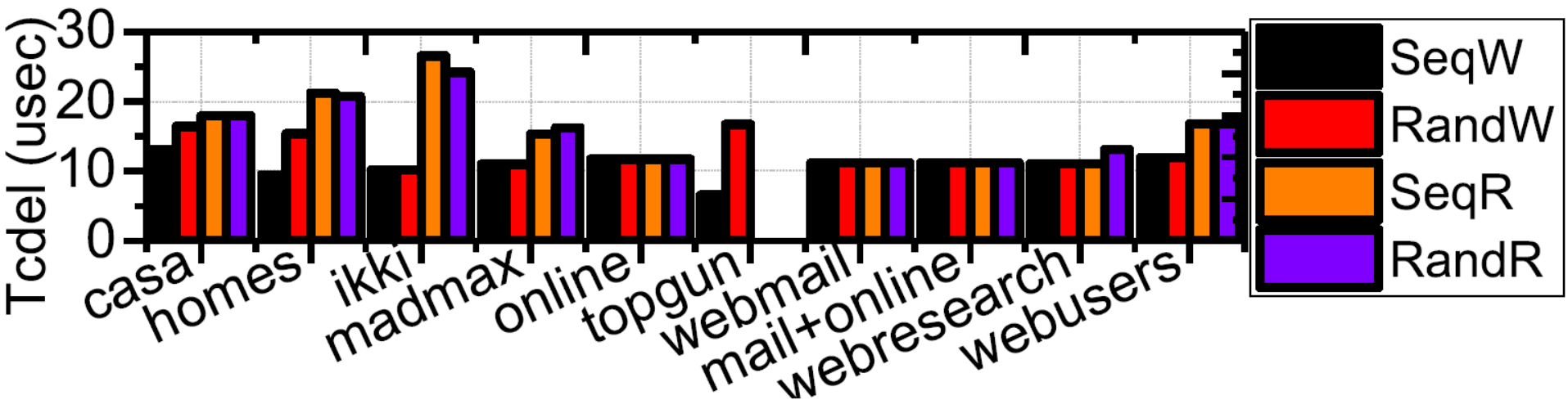
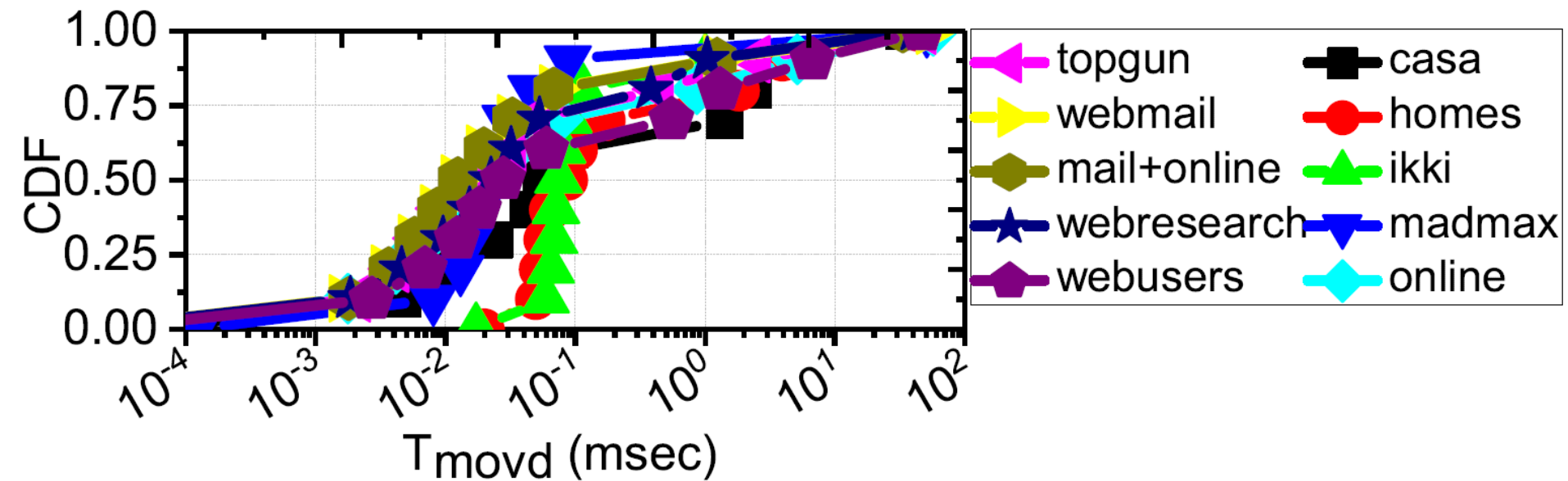
## Hardware/Software Co-Evaluation for Large-Scale I/O Workload Reconstruction

**Miryeong Kwon**, Jie Zhang, Gyuyoung Park,  
Wonil Choi, David Donofrio, John Shalf,  
Mahmut Kandemir, and Myoungsoo Jung



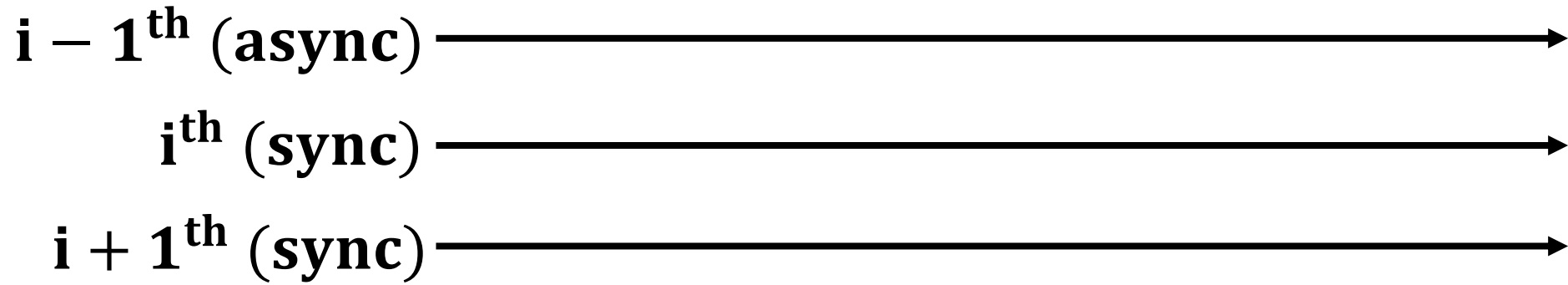
**CAMEL**ab.org  
Computer Architecture and  
Memory Systems  
Laboratory

# CDF distribution of $T_{movd}$ , $T_{cdel}$



# Block Request Timing

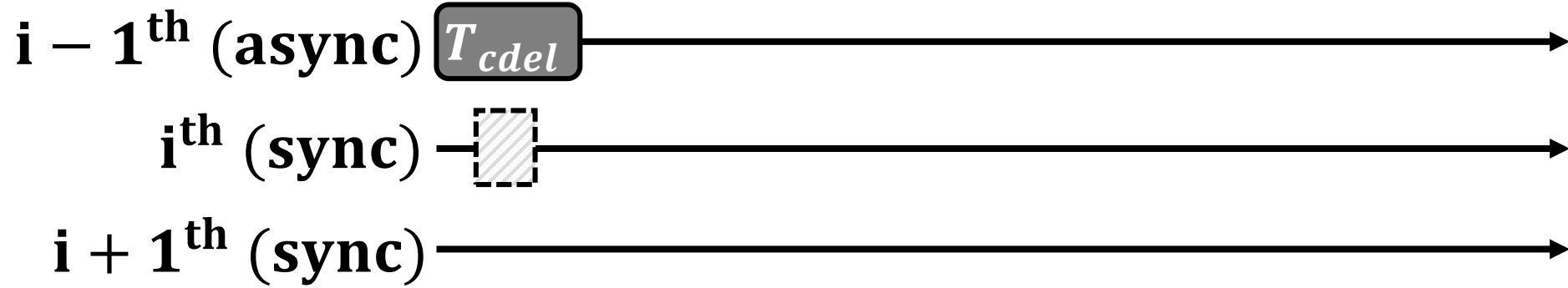
---





# Block Request Timing

---



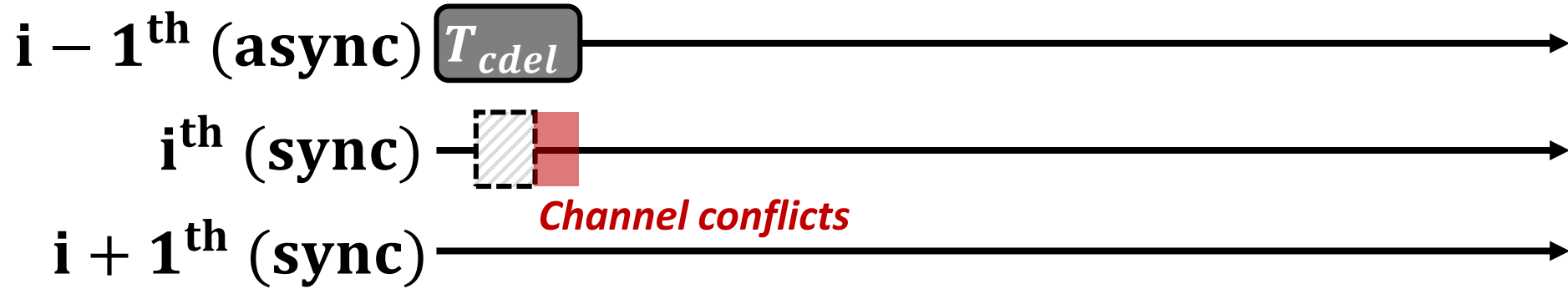
 : Computation cycles by user/kernel

 : Delay due to storage interface (channel) data movement



# Block Request Timing

---



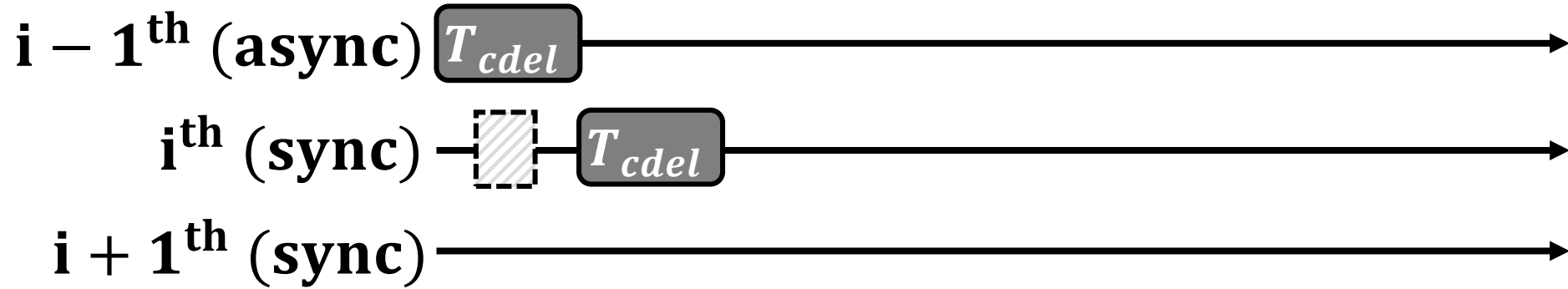
 : Computation cycles by user/kernel

 : Delay due to storage interface (channel) data movement



# Block Request Timing

---



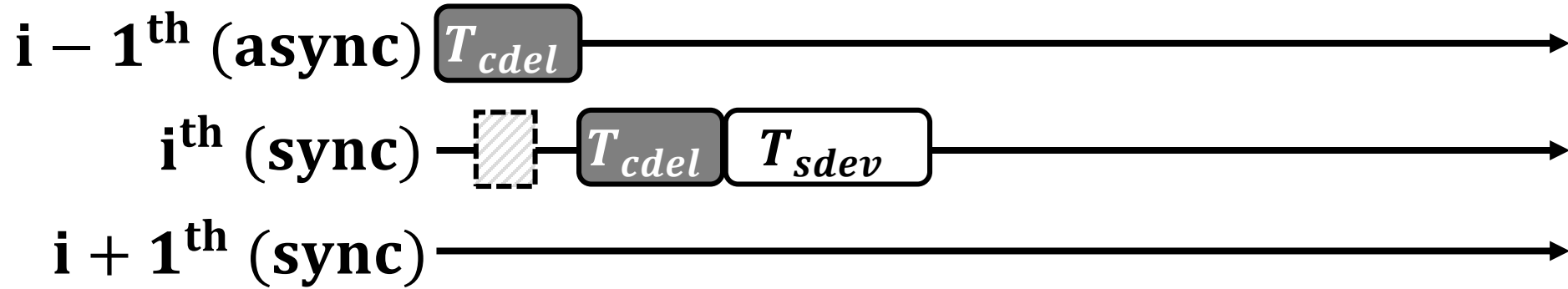
 : Computation cycles by user/kernel

 : Delay due to storage interface (channel) data movement



# Block Request Timing

---



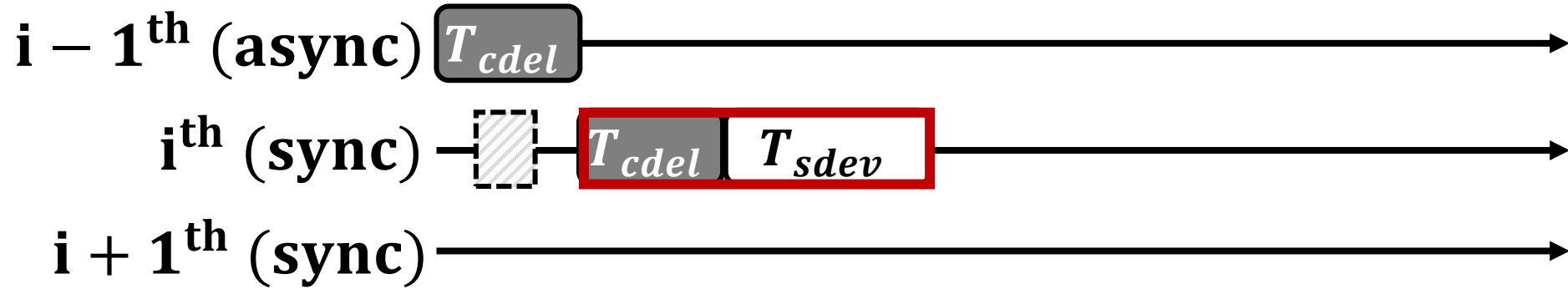
 : Computation cycles by user/kernel

 : Delay due to storage interface (channel) data movement

 : Actual device time taken by the storage to service the IO



# Block Request Timing



 : Computation cycles by user/kernel

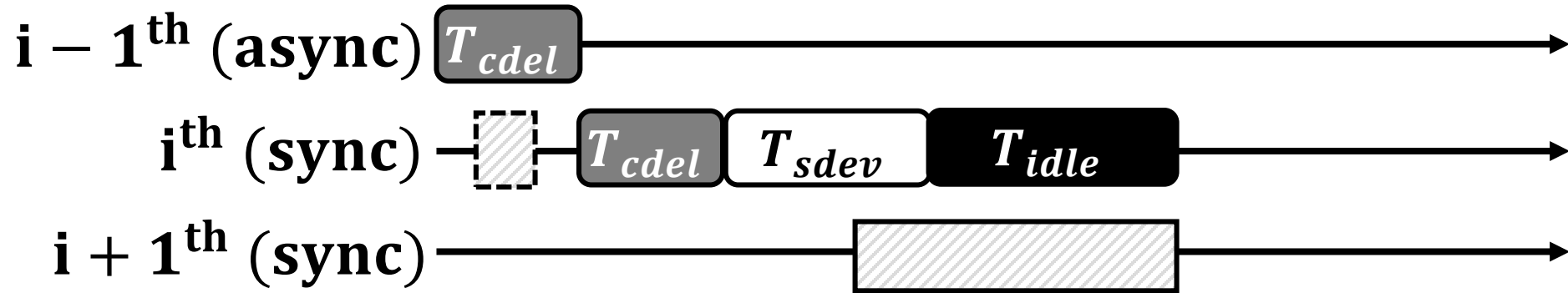
 : Delay due to storage interface (channel) data movement

 : Actual device time taken by the storage to service the IO

*I/O subsystem latency ( $T_{slat}$ )*



# Block Request Timing



 : Computation cycles by user/kernel

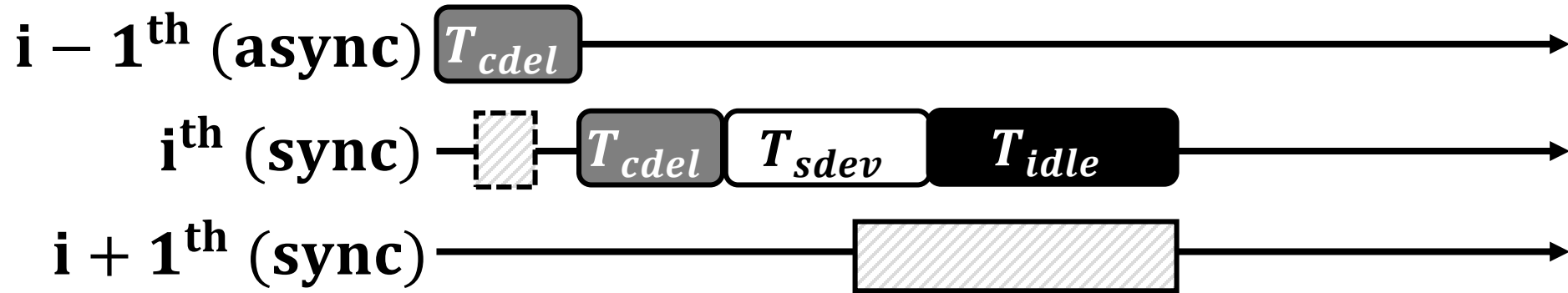
  $T_{c\text{del}}$  : Delay due to storage interface (channel) data movement

  $T_{s\text{dev}}$  : Actual device time taken by the storage to service the IO

  $T_{\text{idle}}$  : user/application does nothing



# Block Request Timing



 : Computation cycles by user/kernel

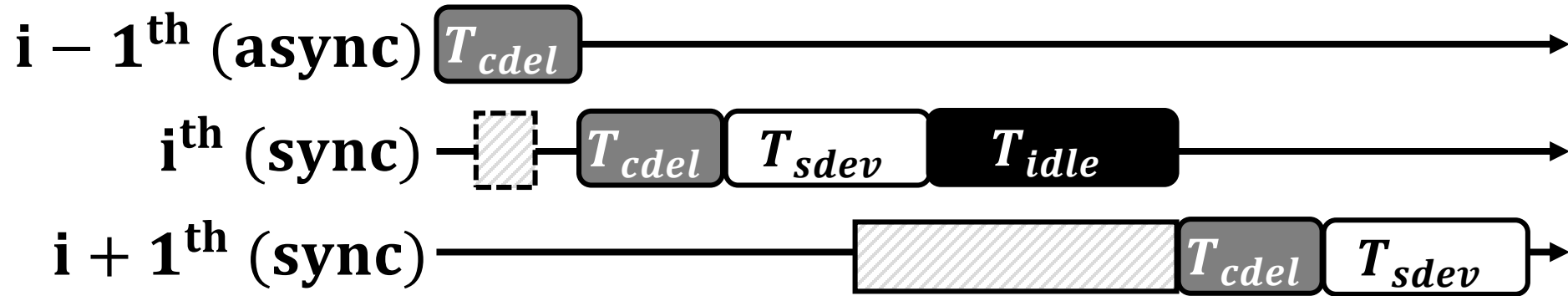
  $T_{cdel}$  : Delay due to storage interface (channel) data movement

  $T_{sdev}$  : Actual device time taken by the storage to service the IO

  $T_{idle}$  : user/application does nothing



# Block Request Timing



 : Computation cycles by user/kernel

 : Delay due to storage interface (channel) data movement

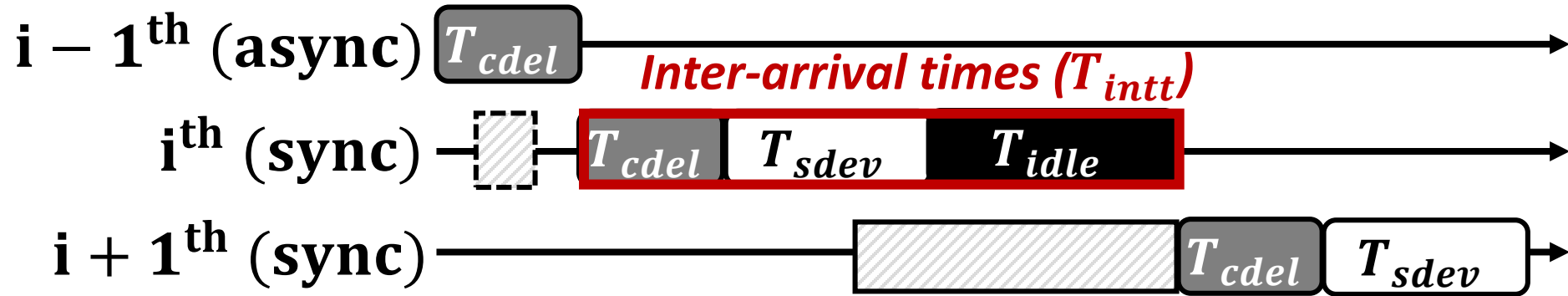
 : Actual device time taken by the storage to service the IO

 : user/application does nothing





# Block Request Timing



 : Computation cycles by user/kernel

 : Delay due to storage interface (channel) data movement

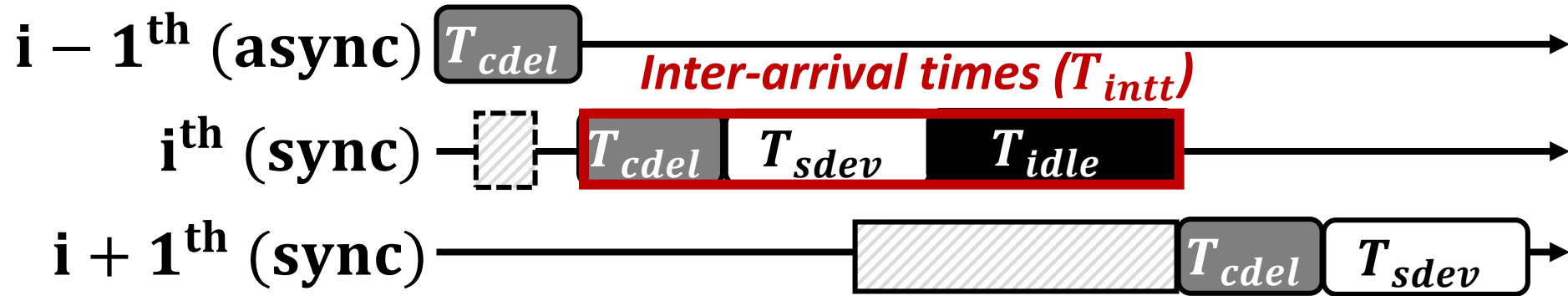
 : Actual device time taken by the storage to service the IO

 : user/application does nothing

Inter-arrival times can be decomposed  
into  $T_{cdel}$ ,  $T_{sdev}$ , and  $T_{idle}$



# Block Request Timing



 : Computation cycles by user/kernel

 : Delay due to storage interface (channel) data movement

 : Actual device time taken by the storage to service the IO

 : user/application does nothing

**Why inter-arrival times decomposition  
is important to infer runtime context?**



# Overview of TraceTracker

---



# Overview of TraceTracker

---

*Pre-software  
evaluation*



# Overview of TraceTracker

---

*Pre-software  
evaluation*



*Hardware  
evaluation*



# Overview of TraceTracker

---

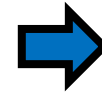


# Overview of TraceTracker

*Pre-software  
evaluation*

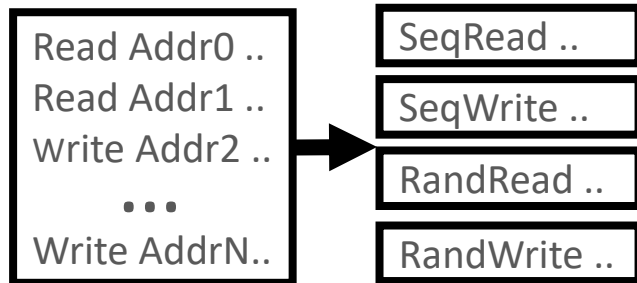


*Hardware  
evaluation*



*Post-software  
evaluation*

Old traces



*Request Classification*

Each  $T_{cdel}$  and  $T_{sdev}$  have similar values  
for same operation type and request size

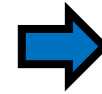


# Overview of TraceTracker

*Pre-software  
evaluation*

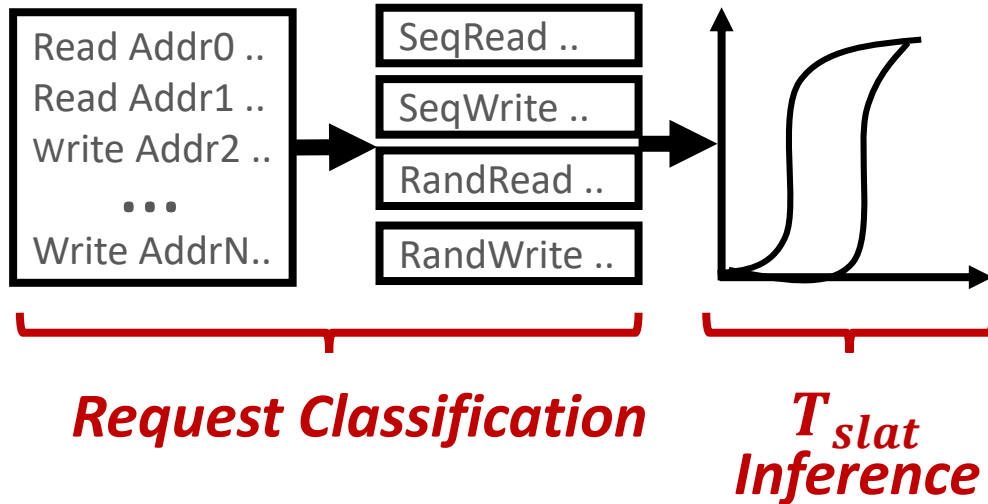


*Hardware  
evaluation*



*Post-software  
evaluation*

Old traces



Infer the  $T_{slat}$  ( $\therefore T_{cdel}$  and  $T_{sdev}$ )  
and calculate  $T_{idle}(=T_{intt}-T_{slat})$





# Overview of TraceTracker

*Pre-software  
evaluation*



*Hardware  
evaluation*

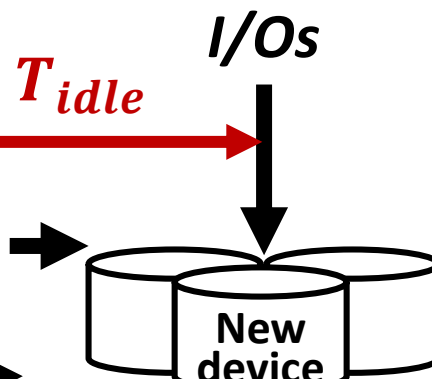
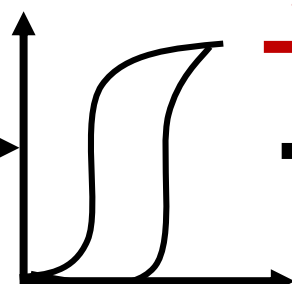


*Post-software  
evaluation*

Old traces

Read Addr0 ..  
Read Addr1 ..  
Write Addr2 ..  
...  
Write AddrN..

SeqRead ..  
SeqWrite ..  
RandRead ..  
RandWrite ..



*Request Classification*

*$T_{slat}$   
Inference*

*Replay  
old traces*

Replay old traces' I/Os with inferred  $T_{idle}$   
while remastering  $T_{slat}$  from new system



# Overview of TraceTracker

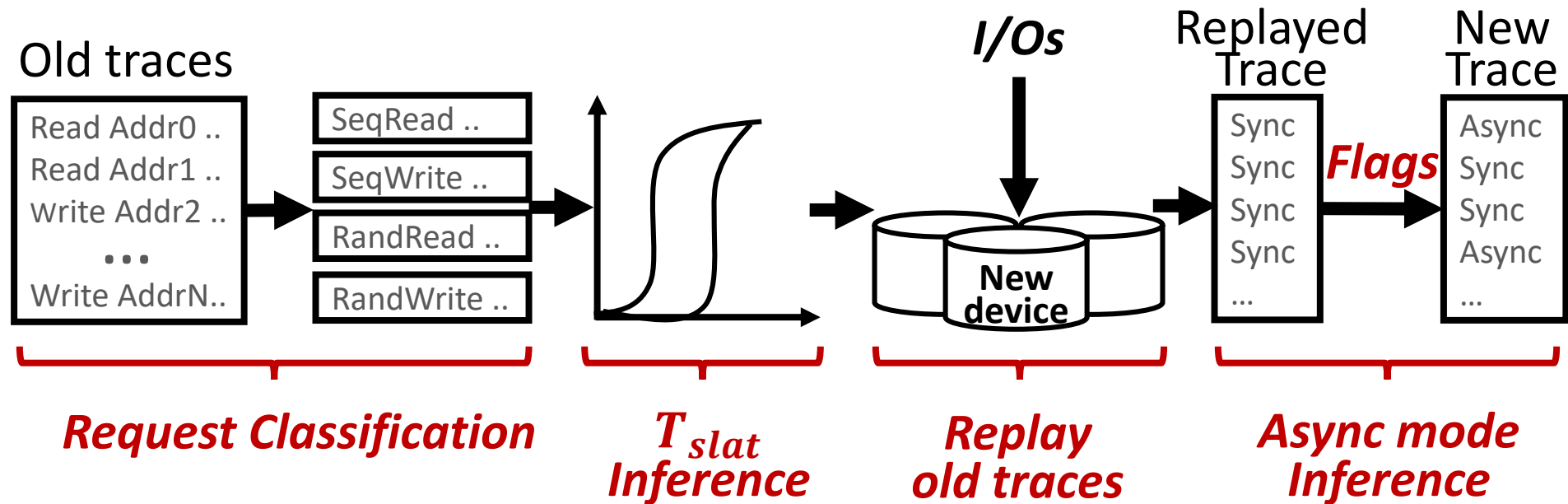
*Pre-software  
evaluation*



*Hardware  
evaluation*



*Post-software  
evaluation*



Infer I/O mode (sync or async) from old traces and revise replayed trace

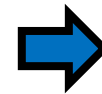


# Overview of TraceTracker

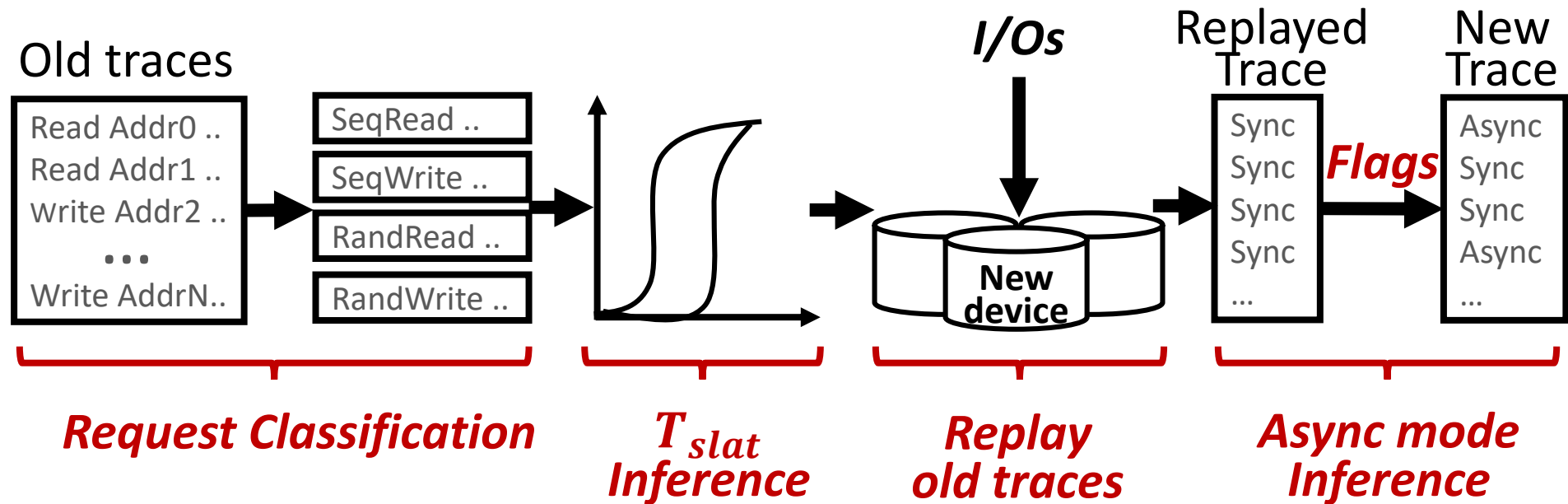
*Pre-software  
evaluation*



*Hardware  
evaluation*



*Post-software  
evaluation*



- 1) Accurate device latency
- 2) Traces include runtime contexts

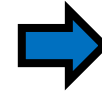


# Overview of TraceTracker

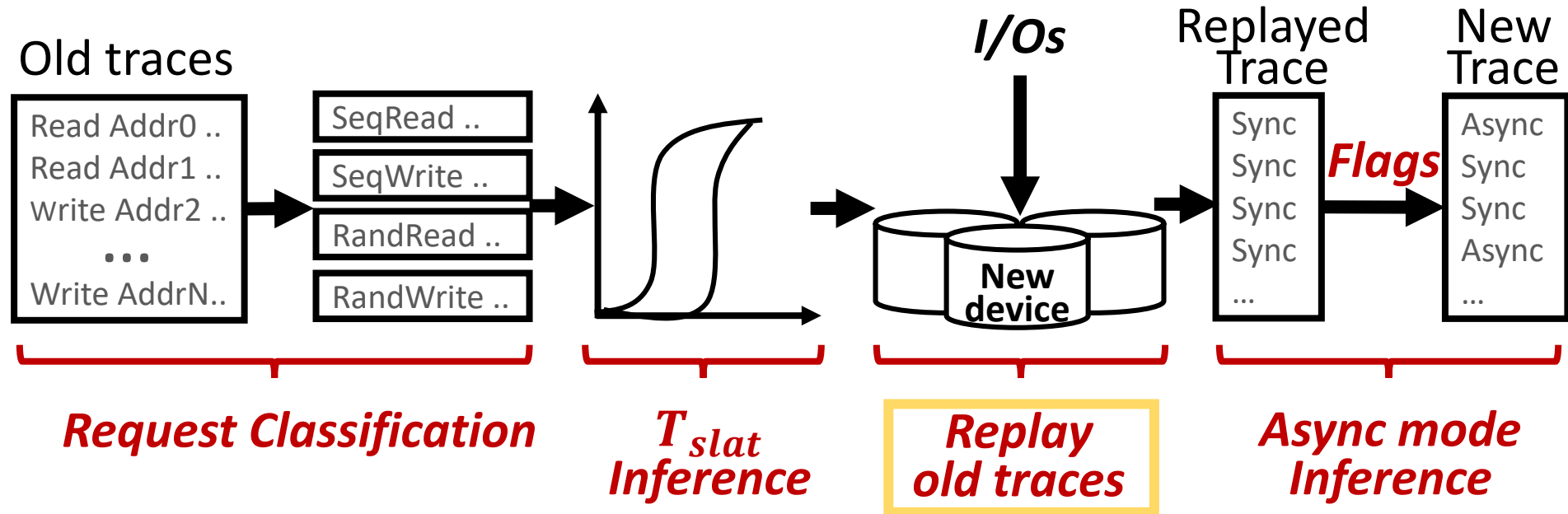
*Pre-software  
evaluation*



*Hardware  
evaluation*



*Post-software  
evaluation*



- 1) Accurate device latency
- 2) Traces include runtime contexts

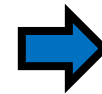


# Overview of TraceTracker

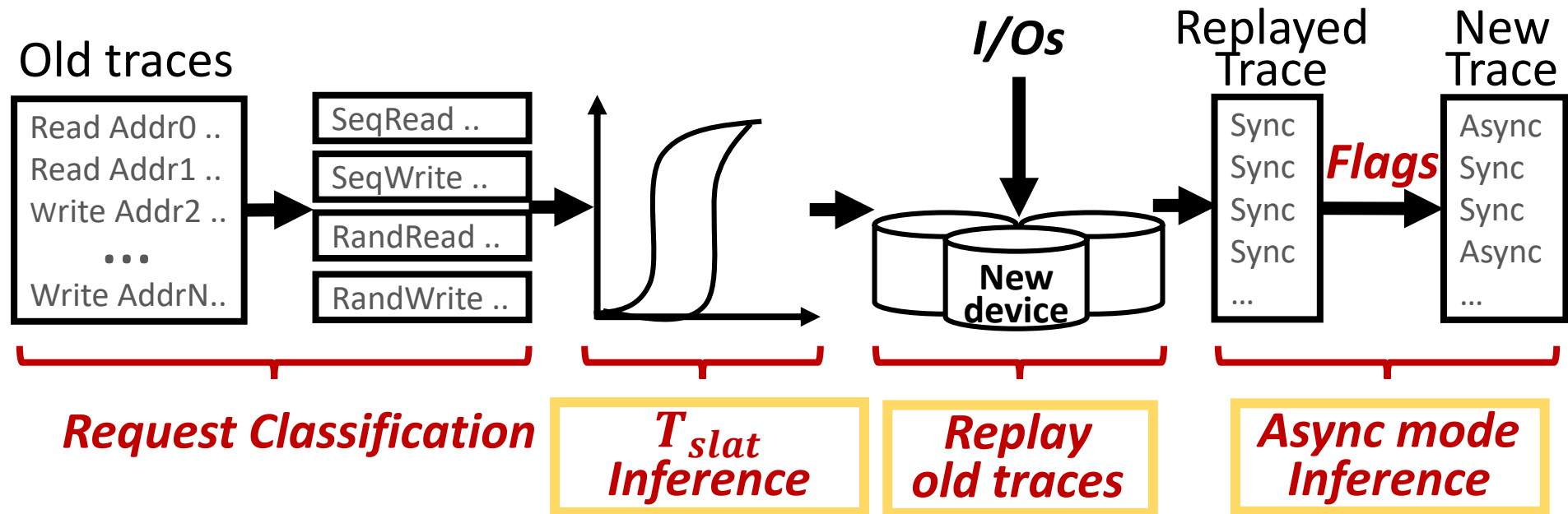
*Pre-software  
evaluation*



*Hardware  
evaluation*



*Post-software  
evaluation*



- 1) Accurate device latency
- 2) Traces include runtime contexts



# Request Classification

---

## Key idea:

- 1) Each request group has similar  $T_{slat}$
- 2)  $T_{sdev}$  can be expressed as linear model



# Request Classification

---

## Key idea:

- 1) Each request group has similar  $T_{slat}$
- 2)  $T_{sdev}$  can be expressed as linear model

$$T_{slat} = T_{cdel} + T_{sdev}$$

## Sequential/Random Read

$$T_{slat} = T_{cdel\_read} + Coeff_{read} \times REQSIZE + T_{movd}$$

## Sequential/Random Write

$$T_{slat} = T_{cdel\_write} + Coeff_{write} \times REQSIZE + T_{movd}$$



# Request Classification

---

## Key idea:

- 1) Each request group has similar  $T_{slat}$
- 2)  $T_{sdev}$  can be expressed as linear model

$$T_{slat} = T_{cdel} + T_{sdev}$$

## Sequential/Random Read

$$T_{slat} = T_{cdel\_read} + Coeff_{read} \times REQSIZE + T_{movd}$$

## Sequential/Random Write

$$T_{slat} = T_{cdel\_write} + Coeff_{write} \times REQSIZE + T_{movd}$$





# Request Classification

---

## Sequential/Random Read

$$T_{slat} = T_{cdel\_read} + Coeff_{read} \times REQSIZE + T_{movd}$$



# Request Classification

---

## Sequential/Random Read

$$T_{slat} = T_{cdel\_read} + Coeff_{read} \times REQSIZE + T_{movd}$$

$$T_{slat}1 = T_{cdel\_read} + Coeff_{read} \times REQSIZE1 + T_{movd}$$

$$T_{slat}2 = T_{cdel\_read} + Coeff_{read} \times REQSIZE2 + T_{movd}$$



# Request Classification

---

## Sequential/Random Read

$$T_{slat} = T_{cdel\_read} + Coeff_{read} \times REQSIZE + T_{movd}$$

$$T_{slat}1 = T_{cdel\_read} + Coeff_{read} \times REQSIZE1 + T_{movd}$$

$$T_{slat}2 = T_{cdel\_read} + Coeff_{read} \times REQSIZE2 + T_{movd}$$



# Request Classification

---

## Sequential/**Random** Read

$$T_{slat} = T_{cdel\_read} + Coeff_{read} \times REQSIZE + T_{movd}$$

$$T_{slat}1 = T_{cdel\_read} + Coeff_{read} \times REQSIZE1 + T_{movd}$$

$$T_{slat}2 = T_{cdel\_read} + Coeff_{read} \times REQSIZE2 + T_{movd}$$

$$T_{slat}1 - T_{slat}2 = Coeff_{read} \times Diff(SIZE1 - SIZE2)$$



# Request Classification

---

## Sequential/Random Read

$$T_{slat} = T_{cdel\_read} + Coeff_{read} \times REQSIZE + T_{movd}$$

$$T_{slat}1 = T_{cdel\_read} + Coeff_{read} \times REQSIZE1 + T_{movd}$$

$$T_{slat}2 = T_{cdel\_read} + Coeff_{read} \times REQSIZE2 + T_{movd}$$

$$T_{slat}1 - T_{slat}2 = Coeff_{read} \times Diff(SIZE1 - SIZE2)$$

*Inferred from max slope of CDF*

*Known info.*



# Request Classification

## Sequential/Random Read

$$T_{slat} = T_{cdel\_read} + Coeff_{read} \times REQSIZE + T_{movd}$$

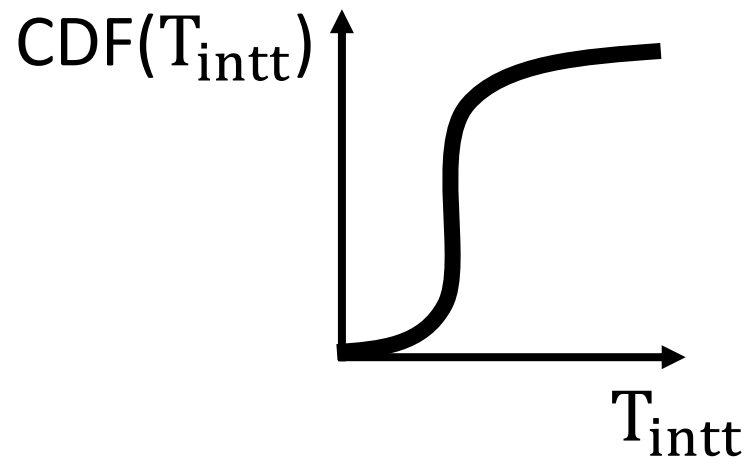
$$T_{slat}1 = T_{cdel\_read} + Coeff_{read} \times REQSIZE1 + T_{movd}$$

$$T_{slat}2 = T_{cdel\_read} + Coeff_{read} \times REQSIZE2 + T_{movd}$$

$$T_{slat}1 - T_{slat}2 = Coeff_{read} \times Diff(SIZE1 - SIZE2)$$

*Inferred from max slope of CDF*

*Known info.*



# Request Classification

## Sequential/Random Read

$$T_{slat} = T_{cdel\_read} + Coeff_{read} \times REQSIZE + T_{movd}$$

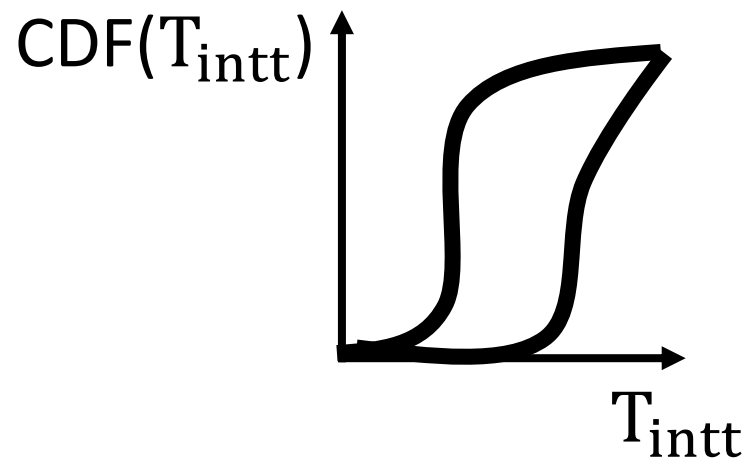
$$T_{slat}1 = T_{cdel\_read} + Coeff_{read} \times REQSIZE1 + T_{movd}$$

$$T_{slat}2 = T_{cdel\_read} + Coeff_{read} \times REQSIZE2 + T_{movd}$$

$$T_{slat}1 - T_{slat}2 = Coeff_{read} \times Diff(SIZE1 - SIZE2)$$

*Inferred from max slope of CDF*

*Known info.*



# Request Classification

## Sequential/Random Read

$$T_{slat} = T_{cdel\_read} + Coeff_{read} \times REQSIZE + T_{movd}$$

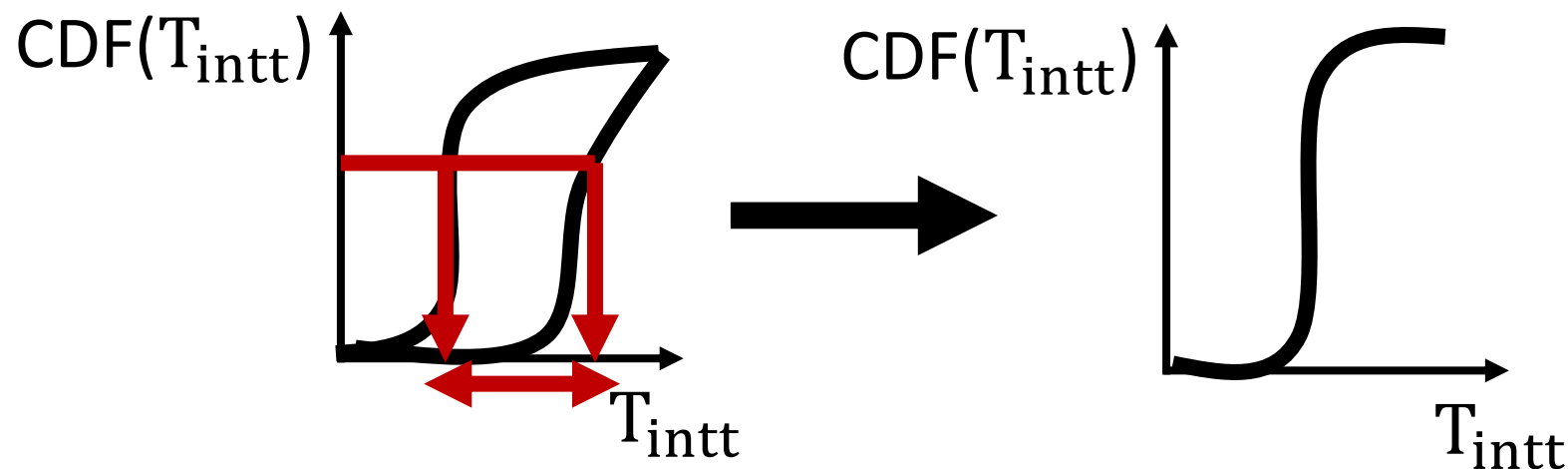
$$T_{slat}1 = T_{cdel\_read} + Coeff_{read} \times REQSIZE1 + T_{movd}$$

$$T_{slat}2 = T_{cdel\_read} + Coeff_{read} \times REQSIZE2 + T_{movd}$$

$$T_{slat}1 - T_{slat}2 = Coeff_{read} \times Diff(SIZE1 - SIZE2)$$

*Inferred from max slope of CDF*

*Known info.*





# Request Classification

## Sequential/Random Read

$$T_{slat} = T_{cdel\_read} + Coeff_{read} \times REQSIZE + T_{movd}$$

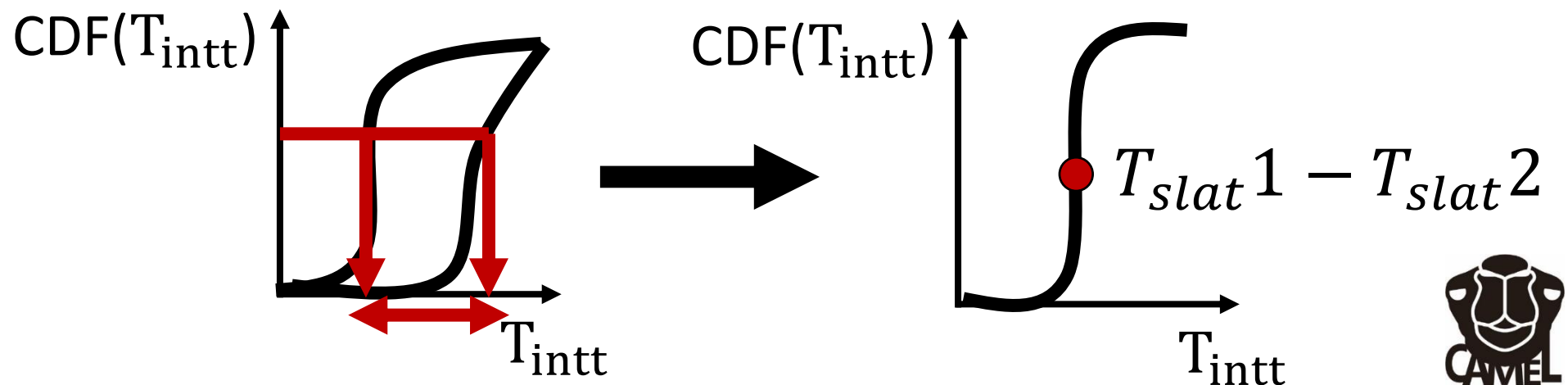
$$T_{slat}1 = T_{cdel\_read} + Coeff_{read} \times REQSIZE1 + T_{movd}$$

$$T_{slat}2 = T_{cdel\_read} + Coeff_{read} \times REQSIZE2 + T_{movd}$$

$$T_{slat}1 - T_{slat}2 = Coeff_{read} \times Diff(SIZE1 - SIZE2)$$

*Inferred from max slope of CDF*

*Known info.*



# Request Classification

---

## Sequential/Random Read

$$T_{slat} = T_{cdel\_read} + Coeff_{read} \times REQSIZE + T_{movd}$$

$$T_{slat}1 = T_{cdel\_read} + Coeff_{read} \times REQSIZE1 + T_{movd}$$

$$T_{slat}2 = T_{cdel\_read} + Coeff_{read} \times REQSIZE2 + T_{movd}$$

$$T_{slat}1 - T_{slat}2 = Coeff_{read} \times Diff(SIZE1 - SIZE2)$$

*Inferred from max slope of CDF*

*Known info.*

Then, which two CDF will be best?

# Inference Model

---

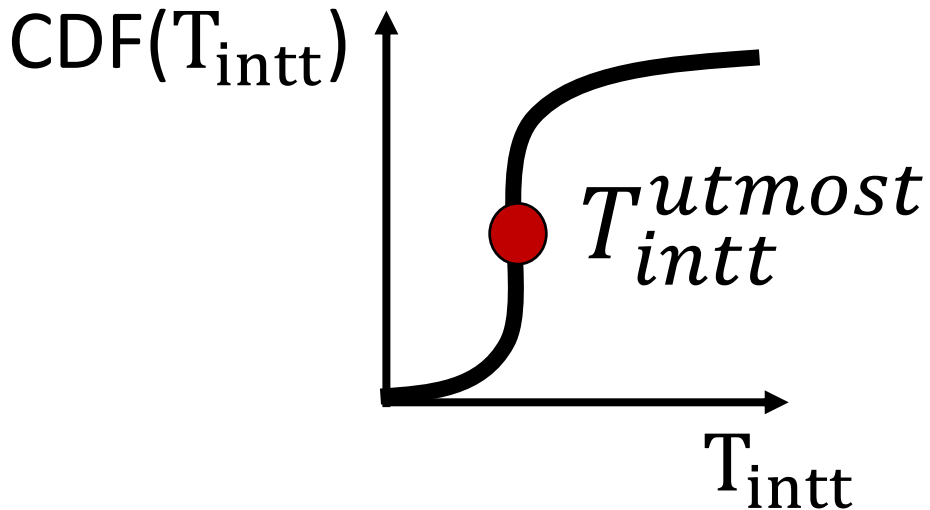
**Key idea:** Select two **steepest** CDF



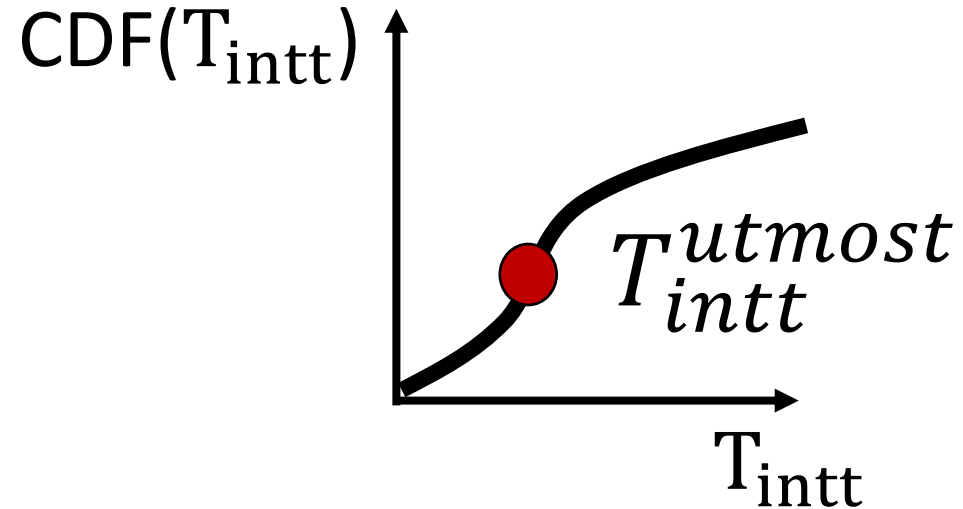
# Inference Model

**Key idea:** Select two **steepest** CDF

$\therefore$  Steep CDF is more accurate to infer  $T_{slat}$



REQ\_SIZE 1



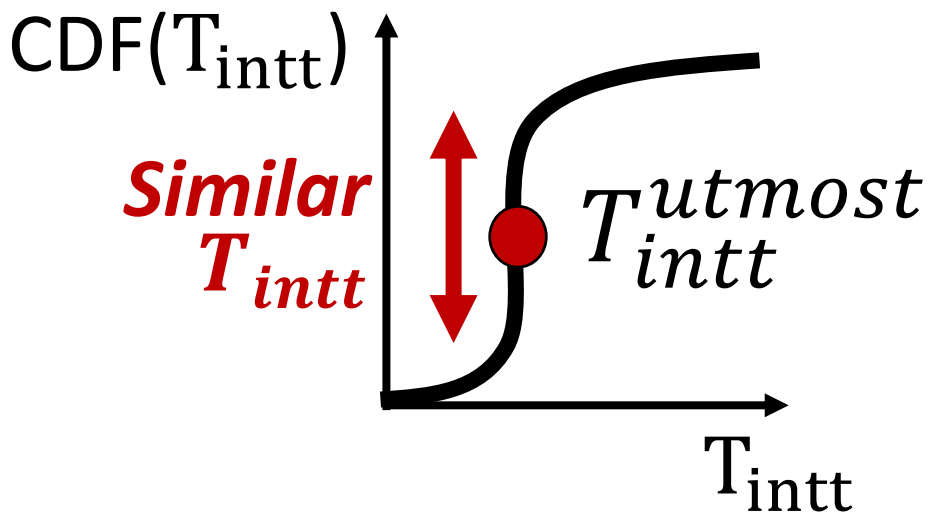
REQ\_SIZE 2



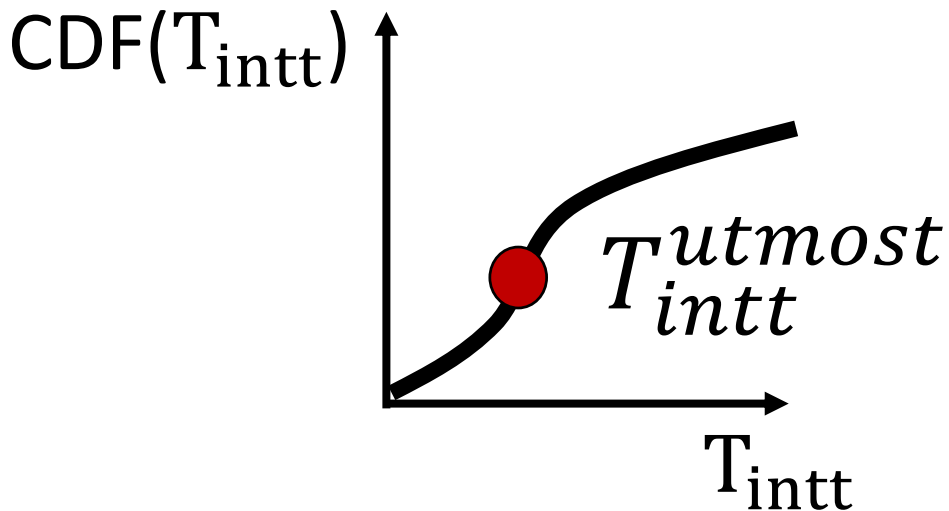
# Inference Model

**Key idea:** Select two **steepest** CDF

$\therefore$  Steep CDF is more accurate to infer  $T_{slat}$



**REQ\_SIZE 1**



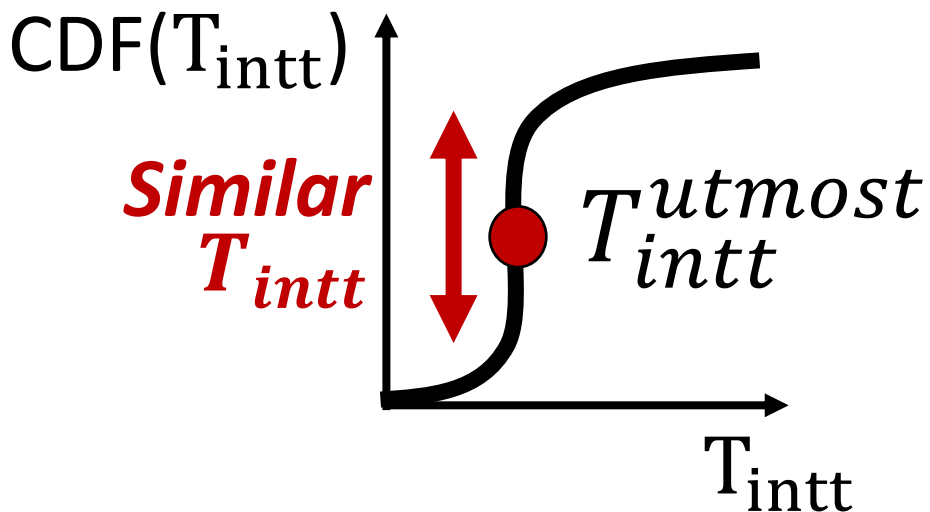
**REQ\_SIZE 2**



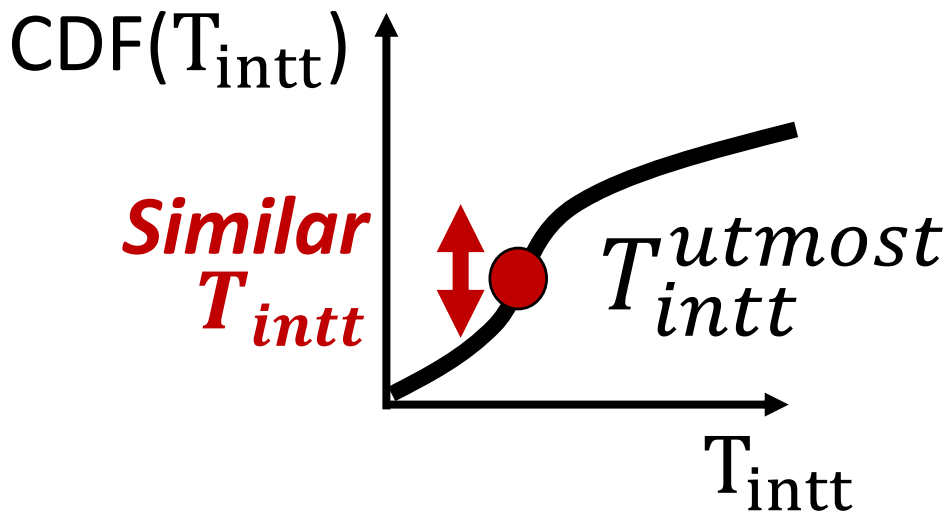
# Inference Model

**Key idea:** Select two **steepest** CDF

$\therefore$  Steep CDF is more accurate to infer  $T_{slat}$



**REQ\_SIZE 1**



**REQ\_SIZE 2**



# Inference Automation

---

**Key idea:** Mathematical analysis;  
Maximum slope of CDF  $\Rightarrow$  Maximum of CDF'



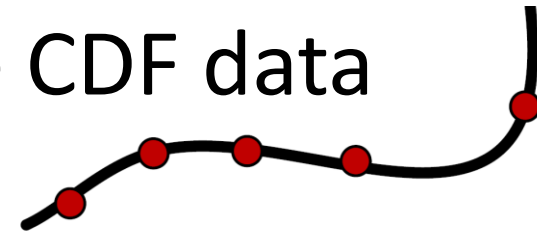
# Inference Automation

---

**Key idea:** Mathematical analysis;  
Maximum slope of CDF  $\Rightarrow$  Maximum of CDF'

**Challenge:** Non-derivative discrete CDF data

**Solution:** Use pchip-interpolation



Pchip interpolation





# Inference Automation

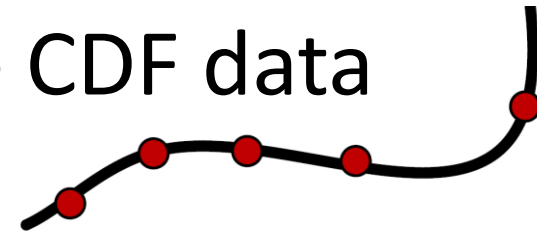
---

**Key idea:** Mathematical analysis;

Maximum slope of CDF  $\Rightarrow$  Maximum of CDF'

**Challenge:** Non-derivative discrete CDF data

**Solution:** Use pchip-interpolation



Pchip interpolation

**Challenge:** High cost of interpolation

**Solution:** Use both PDF and CDF method



# Inference Automation

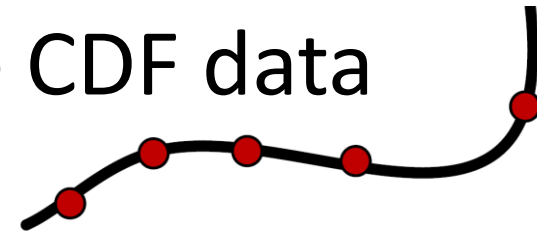
---

**Key idea:** Mathematical analysis;

Maximum slope of CDF  $\Rightarrow$  Maximum of CDF'

**Challenge:** Non-derivative discrete CDF data

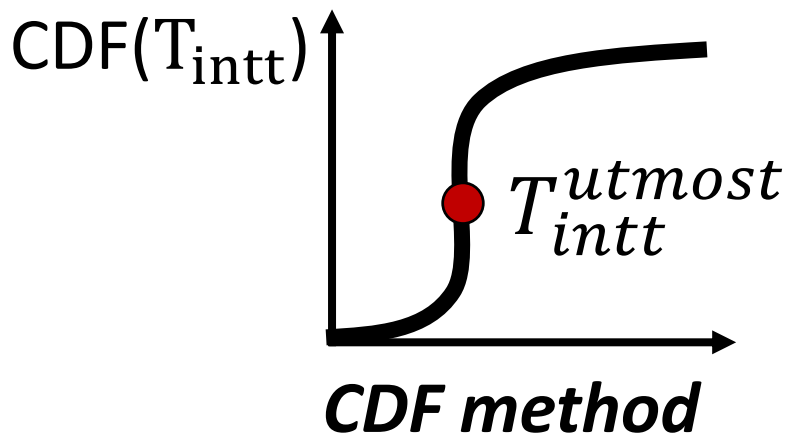
**Solution:** Use pchip-interpolation



Pchip interpolation

**Challenge:** High cost of interpolation

**Solution:** Use both PDF and CDF method



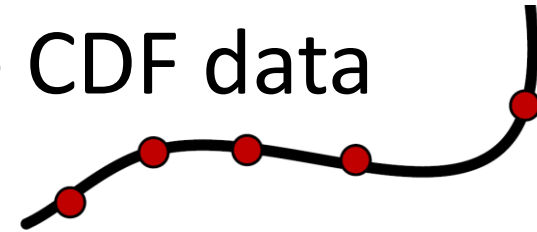
# Inference Automation

**Key idea:** Mathematical analysis;

Maximum slope of CDF => Maximum of CDF'

**Challenge:** Non-derivative discrete CDF data

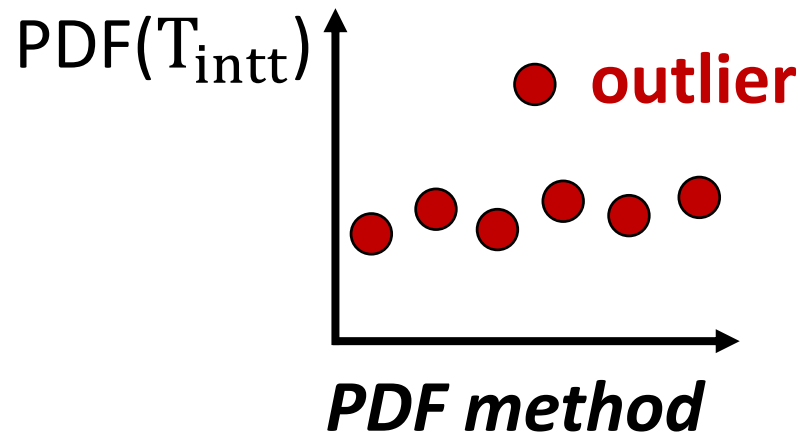
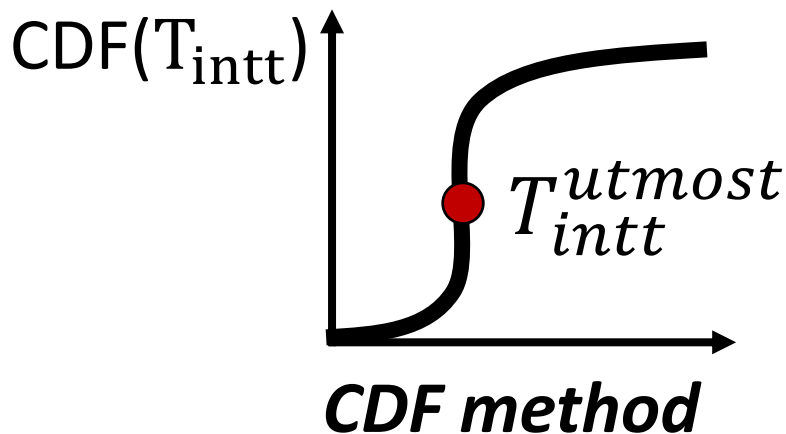
**Solution:** Use pchip-interpolation



Pchip interpolation

**Challenge:** High cost of interpolation

**Solution:** Use both PDF and CDF method



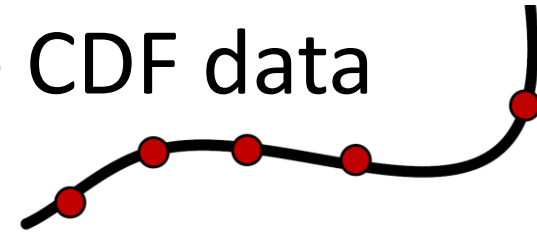
# Inference Automation

**Key idea:** Mathematical analysis;

Maximum slope of CDF => Maximum of CDF'

**Challenge:** Non-derivative discrete CDF data

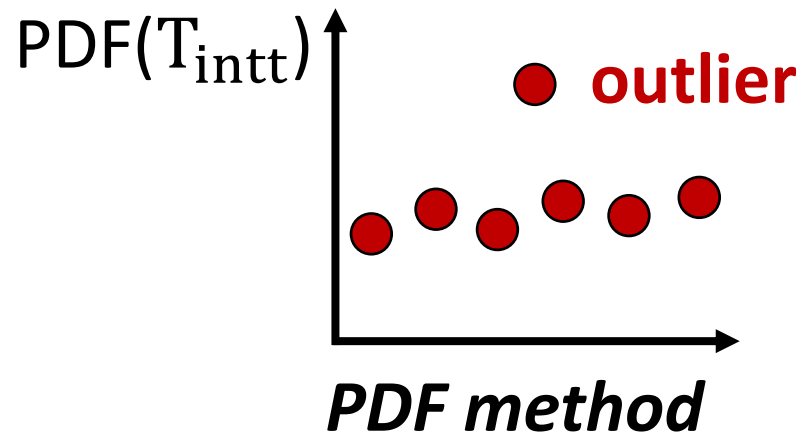
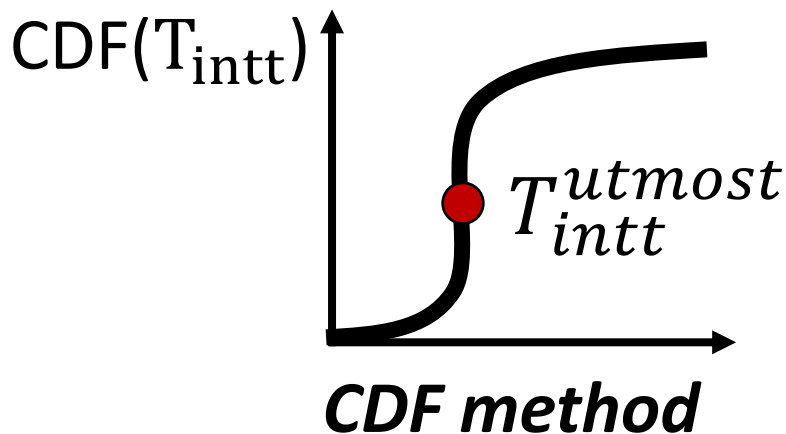
**Solution:** Use pchip-interpolation



Pchip interpolation

**Challenge:** High cost of interpolation

**Solution:** Use both PDF and CDF method



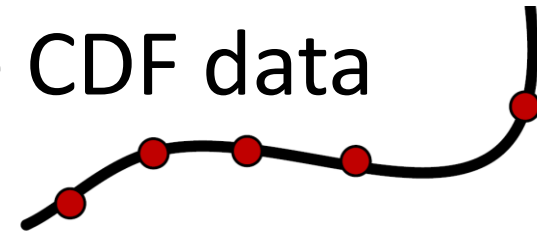
# Inference Automation

**Key idea:** Mathematical analysis;

Maximum slope of CDF => Maximum of CDF'

**Challenge:** Non-derivative discrete CDF data

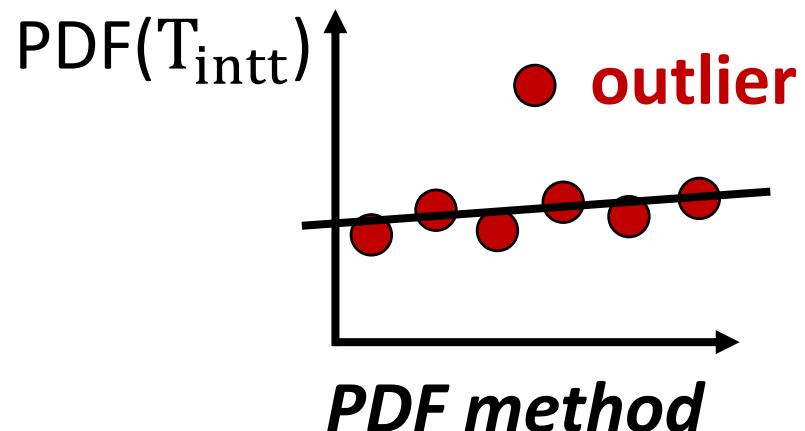
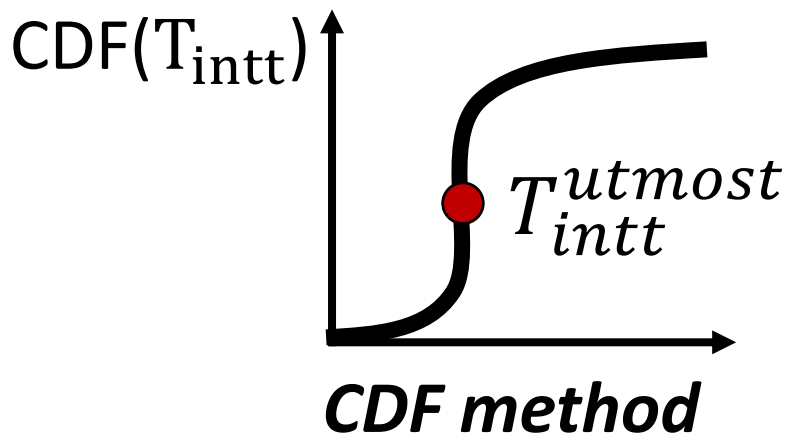
**Solution:** Use pchip-interpolation



Pchip interpolation

**Challenge:** High cost of interpolation

**Solution:** Use both PDF and CDF method



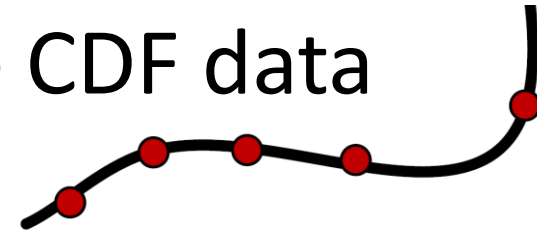
# Inference Automation

**Key idea:** Mathematical analysis;

Maximum slope of CDF => Maximum of CDF'

**Challenge:** Non-derivative discrete CDF data

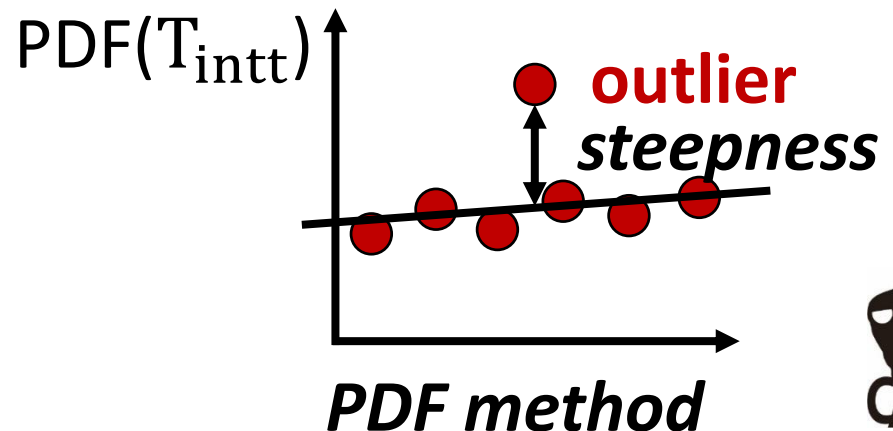
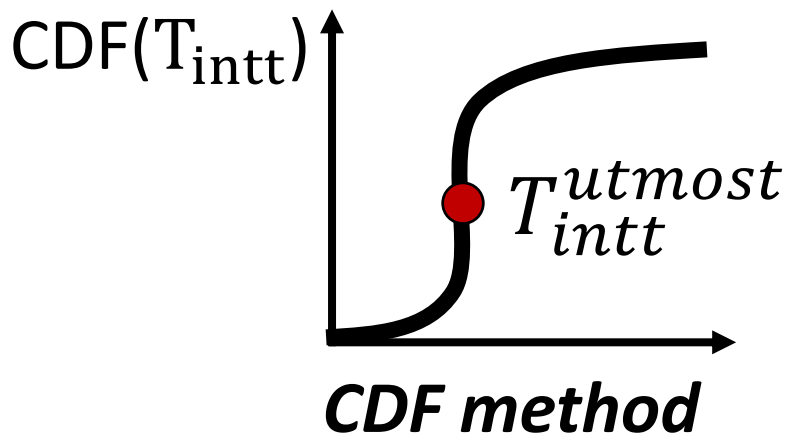
**Solution:** Use pchip-interpolation



Pchip interpolation

**Challenge:** High cost of interpolation

**Solution:** Use both PDF and CDF method



# Inference Automation

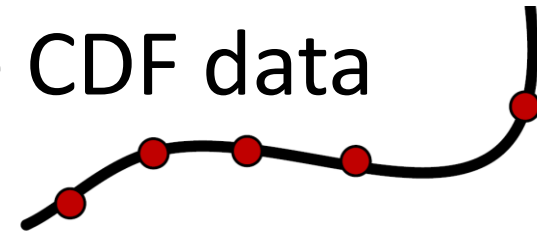
---

**Key idea:** Mathematical analysis;

Maximum slope of CDF => Maximum of CDF'

**Challenge:** Non-derivative discrete CDF data

**Solution:** Use pchip-interpolation



Pchip interpolation

**Challenge:** High cost of interpolation

**Solution:**

	PDF method	CDF method
Accuracy	<	
Efficiency	>	



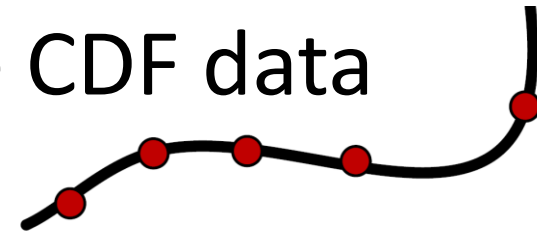
# Inference Automation

**Key idea:** Mathematical analysis;

Maximum slope of CDF  $\Rightarrow$  Maximum of CDF'

**Challenge:** Non-derivative discrete CDF data

**Solution:** Use pchip-interpolation



Pchip interpolation

**Challenge:** High cost of interpolation

**Solution:**

	PDF method	CDF method
Accuracy	<	
Efficiency	>	

*Steepness analysis*





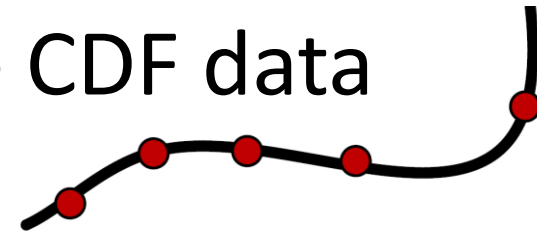
# Inference Automation

**Key idea:** Mathematical analysis;

Maximum slope of CDF  $\Rightarrow$  Maximum of CDF'

**Challenge:** Non-derivative discrete CDF data

**Solution:** Use pchip-interpolation



Pchip interpolation

**Challenge:** High cost of interpolation

**Solution:**

	PDF method	CDF method
Accuracy	<	
Efficiency	>	

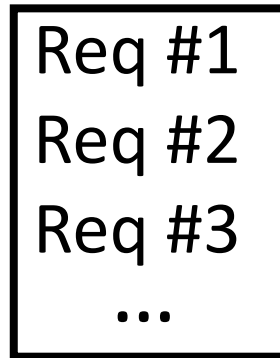
*Find max slope of CDF(diff)*



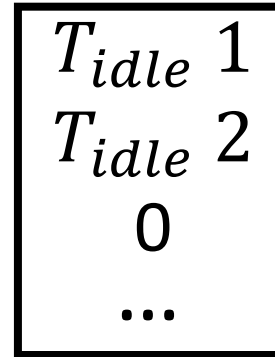
# Verification

---

Old traces



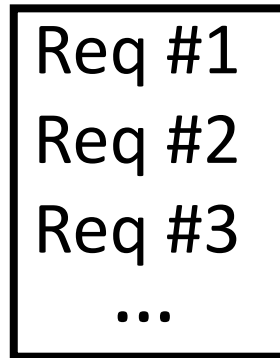
Injected  $T_{idle}$



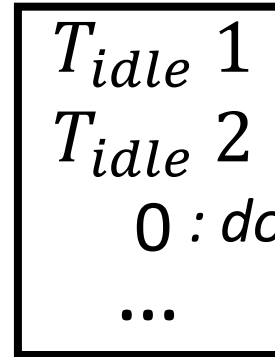
# Verification

---

Old traces



Injected  $T_{idle}$

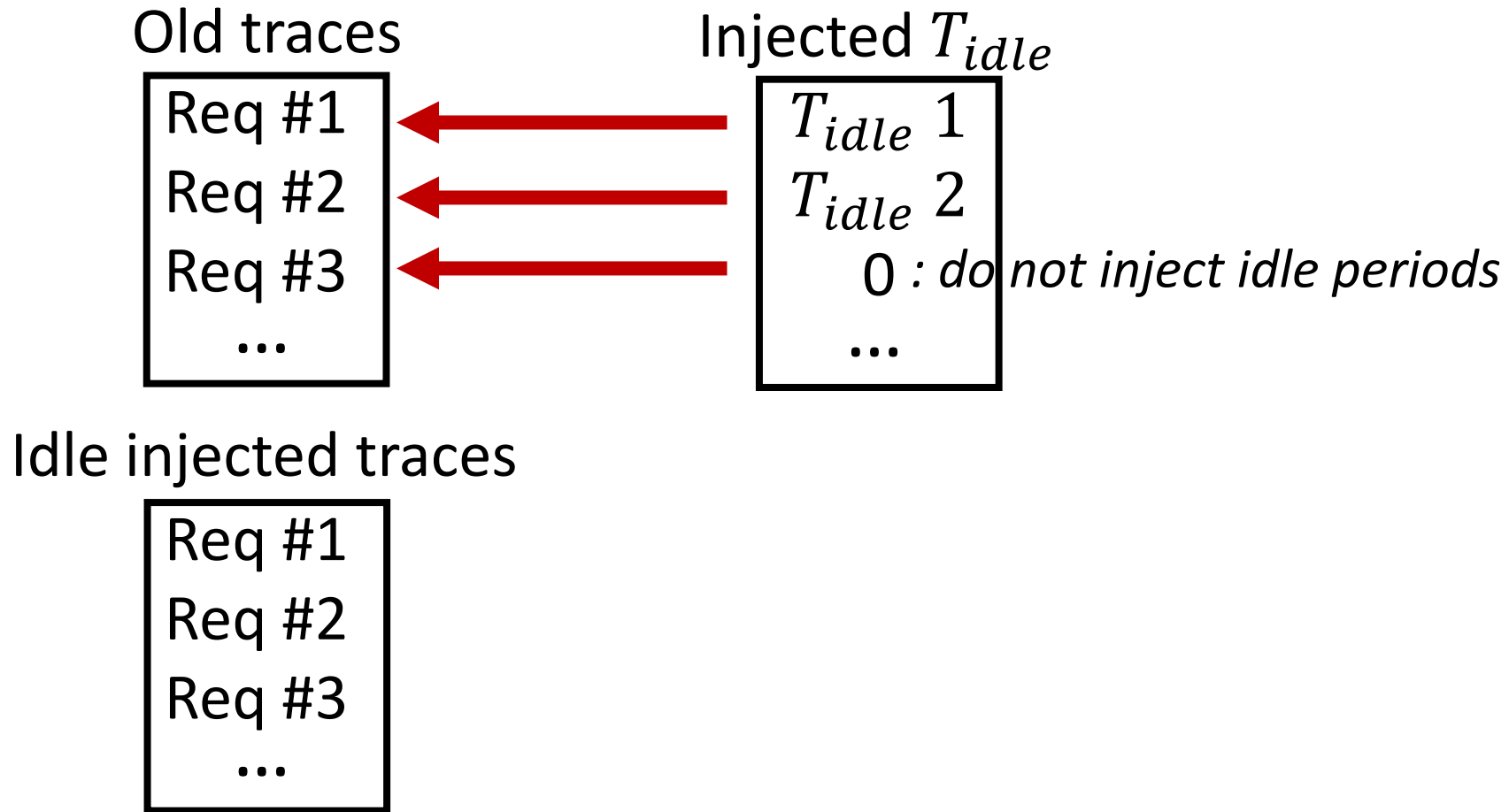


*0 : do not inject idle periods*



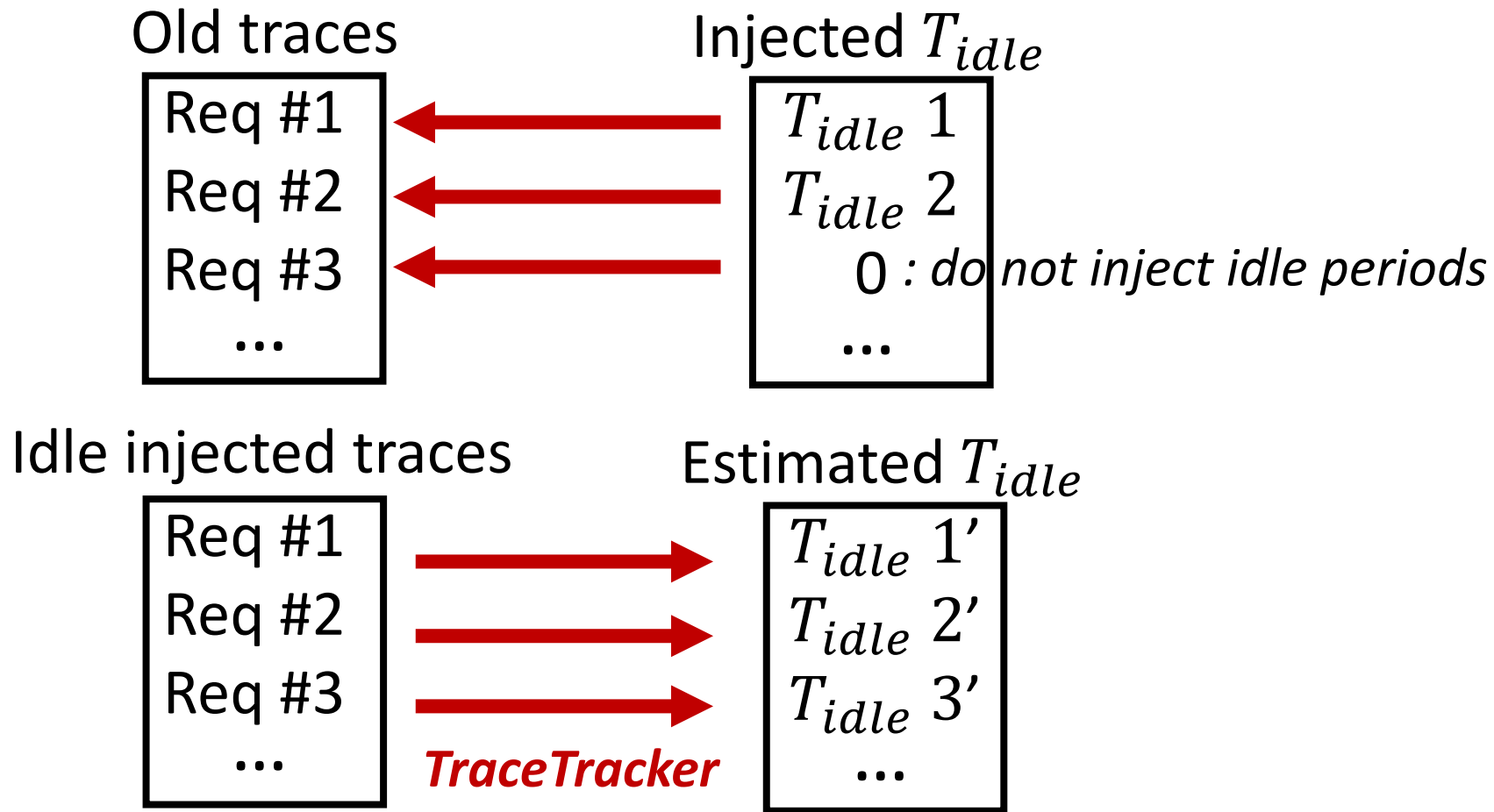
# Verification

---

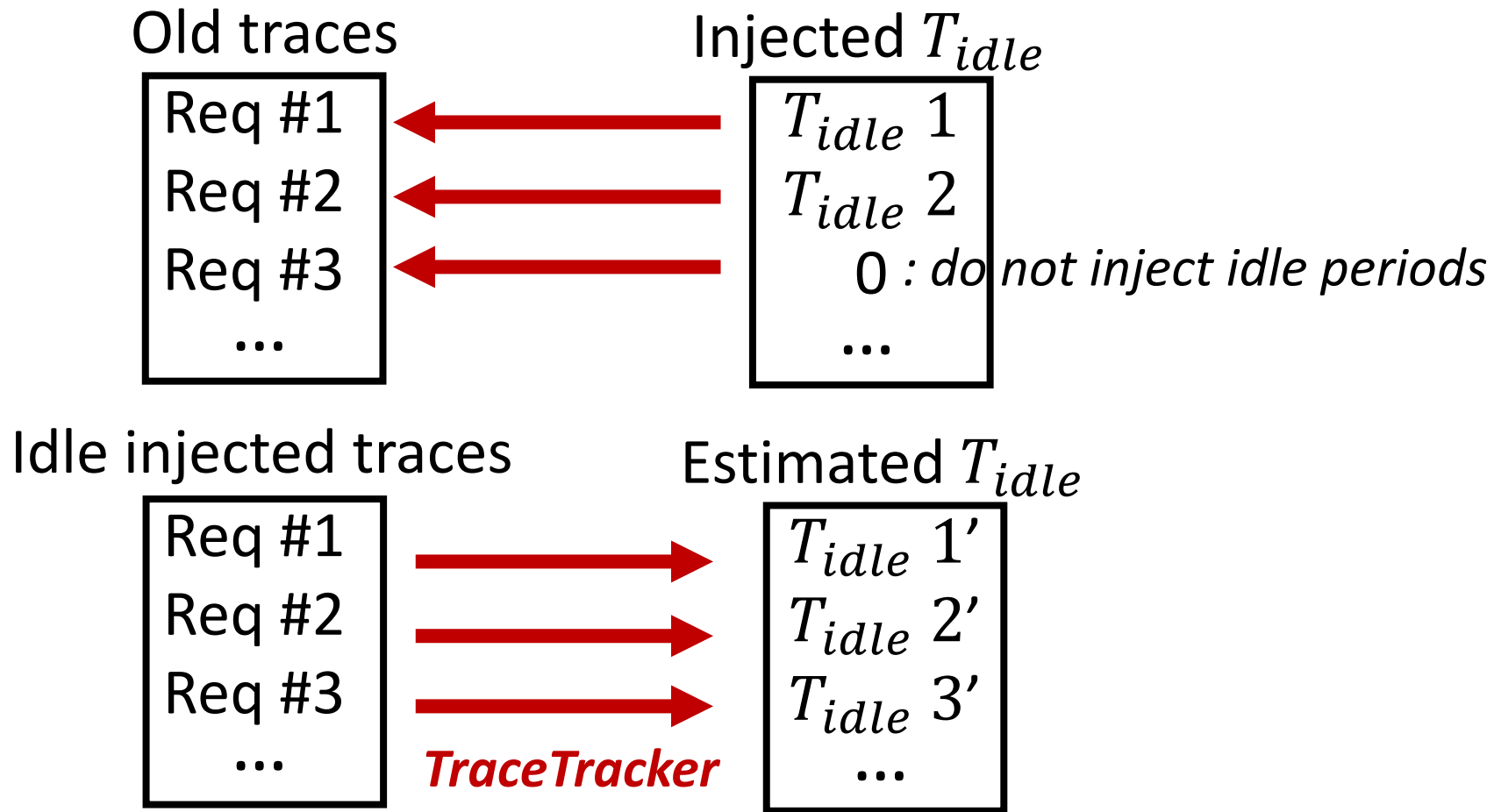


# Verification

---



# Verification



- 1) Number of  $T_{idle}$  occurrence: 99%
- 2) Total periods of  $T_{idle}$ : 96%

