



# ALTERNATE PATH $\mu$ -OP CACHE PREFETCHING

**Sawan Singh**<sup>1</sup>    **Arthur Perais**<sup>2</sup>    **Alexandra Jimborean**<sup>1</sup>  
**Alberto Ros**<sup>1</sup>

<sup>1</sup>Computer Engineering Department  
University of Murcia, Spain

<sup>2</sup>TIMA, Univ. Grenoble Alpes, CNRS,  
Grenoble INP, Grenoble, France



ISCA'51, Session 10A, July 3, 2024

# OVERVIEW

→  $\mu$ -op Cache

## OVERVIEW

→  $\mu$ -op Cache

- Holds recently decoded  $\mu$ -ops

## OVERVIEW

### → $\mu$ -op Cache

- Holds recently decoded  $\mu$ -ops
- First introduced for **energy savings**<sup>1</sup> in x86 which requires complex decoder

<sup>1</sup> Solomon et al. *Micro-operation cache: a power aware frontend for variable instruction length ISA*

## OVERVIEW

### → $\mu$ -op Cache

- Holds recently decoded  $\mu$ -ops
- First introduced for **energy savings**<sup>1</sup> in x86 which requires complex decoder
- Current  $\mu$ -op caches are severely overwhelmed by **server workloads**

<sup>1</sup> Solomon et al. *Micro-operation cache: a power aware frontend for variable instruction length ISA*

## OVERVIEW

### → $\mu$ -op Cache

- Holds recently decoded  $\mu$ -ops
- First introduced for **energy savings**<sup>1</sup> in x86 which requires complex decoder
- Current  $\mu$ -op caches are severely overwhelmed by **server workloads**
  - Only provide **0.87%** improvement over no  $\mu$ -op cache

<sup>1</sup> Solomon et al. *Micro-operation cache: a power aware frontend for variable instruction length ISA*

## OVERVIEW

### → $\mu$ -op Cache

- Holds recently decoded  $\mu$ -ops
- First introduced for **energy savings**<sup>1</sup> in x86 which requires complex decoder
- Current  $\mu$ -op caches are severely overwhelmed by **server workloads**
  - Only provide **0.87%** improvement over no  $\mu$ -op cache
  - Ideal  $\mu$ -op cache can provide **10.82%** improvement

<sup>1</sup> Solomon et al. Micro-operation cache: a power aware frontend for variable instruction length ISA

## OVERVIEW

### → $\mu$ -op Cache

- Holds recently decoded  $\mu$ -ops
- First introduced for **energy savings**<sup>1</sup> in x86 which requires complex decoder
- Current  $\mu$ -op caches are severely overwhelmed by **server workloads**
  - Only provide **0.87%** improvement over no  $\mu$ -op cache
  - Ideal  $\mu$ -op cache can provide **10.82%** improvement

→ We propose **UCP** (*Alternate Path  $\mu$ -op Cache Prefetching*)

<sup>1</sup> Solomon et al. Micro-operation cache: a power aware frontend for variable instruction length ISA



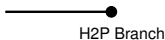
## OVERVIEW

### → $\mu$ -op Cache

- Holds recently decoded  $\mu$ -ops
- First introduced for **energy savings**<sup>1</sup> in x86 which requires complex decoder
- Current  $\mu$ -op caches are severely overwhelmed by **server workloads**
  - Only provide **0.87%** improvement over no  $\mu$ -op cache
  - Ideal  $\mu$ -op cache can provide **10.82%** improvement

### → We propose **UCP** (*Alternate Path $\mu$ -op Cache Prefetching*)

- Identify **hard-to-predict** branches



<sup>1</sup> Solomon et al. *Micro-operation cache: a power aware frontend for variable instruction length ISA*

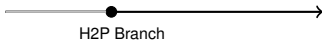
## OVERVIEW

### → $\mu$ -op Cache

- Holds recently decoded  $\mu$ -ops
- First introduced for **energy savings**<sup>1</sup> in x86 which requires complex decoder
- Current  $\mu$ -op caches are severely overwhelmed by **server workloads**
  - Only provide **0.87%** improvement over no  $\mu$ -op cache
  - Ideal  $\mu$ -op cache can provide **10.82%** improvement

### → We propose **UCP** (*Alternate Path $\mu$ -op Cache Prefetching*)

- Identify **hard-to-predict** branches
- Prefetch  $\mu$ -ops from **alternate path** of the hard-to-predict branch



<sup>1</sup> Solomon et al. Micro-operation cache: a power aware frontend for variable instruction length ISA

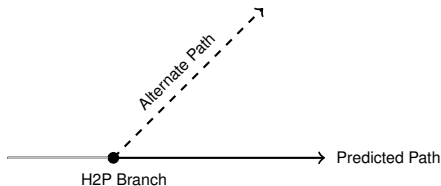
## OVERVIEW

### → $\mu$ -op Cache

- Holds recently decoded  $\mu$ -ops
- First introduced for **energy savings**<sup>1</sup> in x86 which requires complex decoder
- Current  $\mu$ -op caches are severely overwhelmed by **server workloads**
  - Only provide **0.87%** improvement over no  $\mu$ -op cache
  - Ideal  $\mu$ -op cache can provide **10.82%** improvement

### → We propose **UCP** (*Alternate Path $\mu$ -op Cache Prefetching*)

- Identify **hard-to-predict** branches
- Prefetch  $\mu$ -ops from **alternate path** of the hard-to-predict branch



<sup>1</sup> Solomon et al. Micro-operation cache: a power aware frontend for variable instruction length ISA

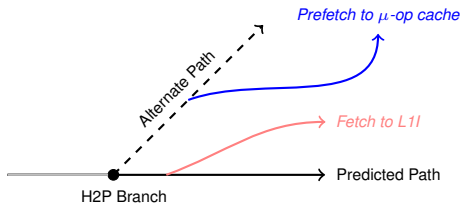
## OVERVIEW

### → $\mu$ -op Cache

- Holds recently decoded  $\mu$ -ops
- First introduced for **energy savings**<sup>1</sup> in x86 which requires complex decoder
- Current  $\mu$ -op caches are severely overwhelmed by **server workloads**
  - Only provide **0.87%** improvement over no  $\mu$ -op cache
  - Ideal  $\mu$ -op cache can provide **10.82%** improvement

### → We propose **UCP** (*Alternate Path $\mu$ -op Cache Prefetching*)

- Identify **hard-to-predict** branches
- Prefetch  $\mu$ -ops from **alternate path** of the hard-to-predict branch



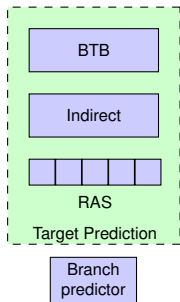
<sup>1</sup> Solomon et al. Micro-operation cache: a power aware frontend for variable instruction length ISA

# OUTLINE

- 1 Overview
- 2 **Background & Motivation**
- 3 UCP
- 4 Methodology & Results
- 5 Conclusions

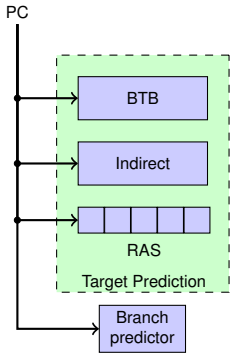
# BACKGROUND & MOTIVATION

## PROCESSOR FRONT-END



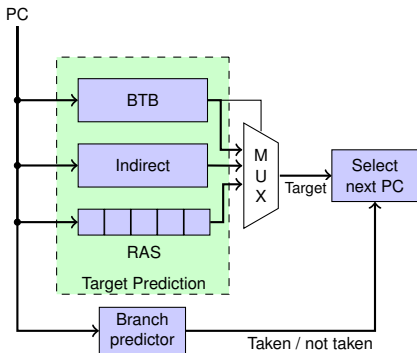
# BACKGROUND & MOTIVATION

## PROCESSOR FRONT-END



# BACKGROUND & MOTIVATION

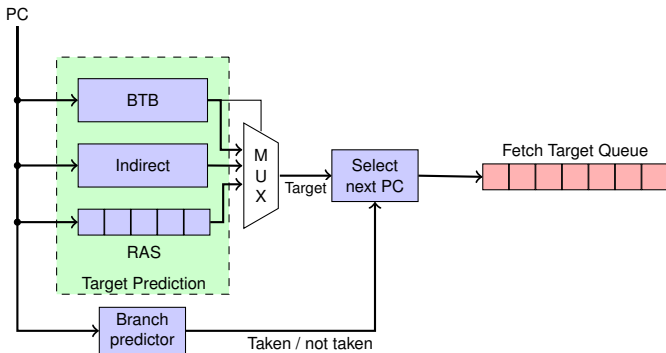
## PROCESSOR FRONT-END





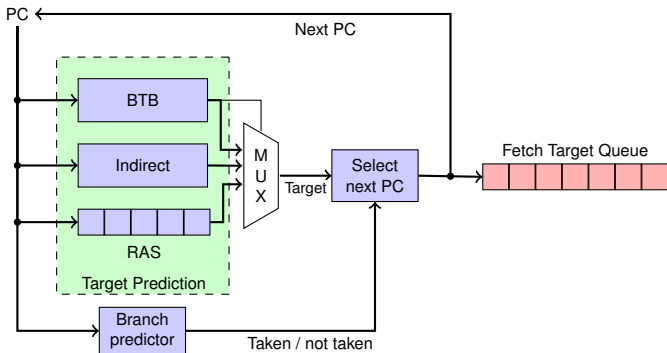
# BACKGROUND & MOTIVATION

## PROCESSOR FRONT-END



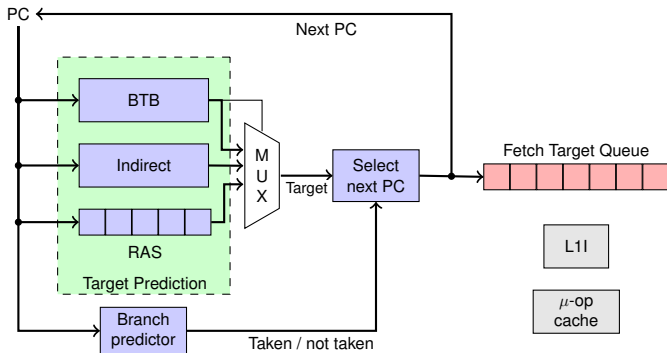
# BACKGROUND & MOTIVATION

## PROCESSOR FRONT-END



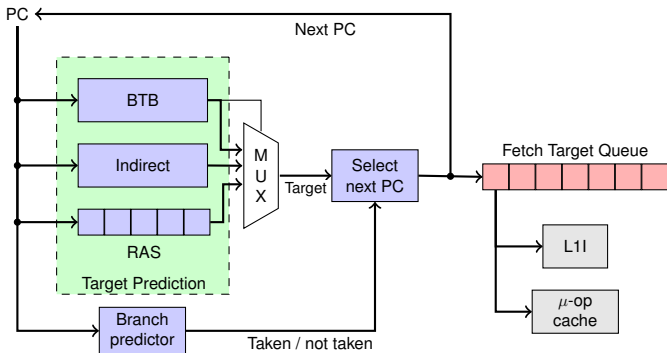
# BACKGROUND & MOTIVATION

## PROCESSOR FRONT-END



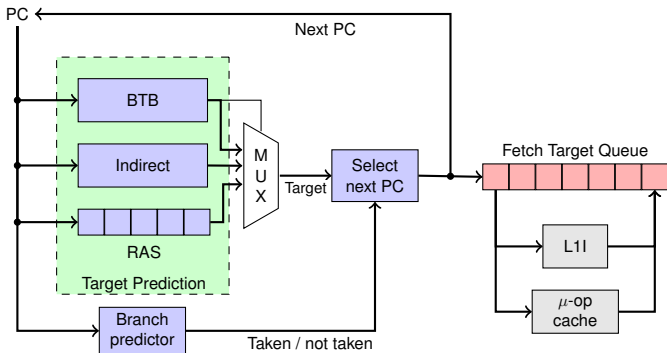
# BACKGROUND & MOTIVATION

## PROCESSOR FRONT-END



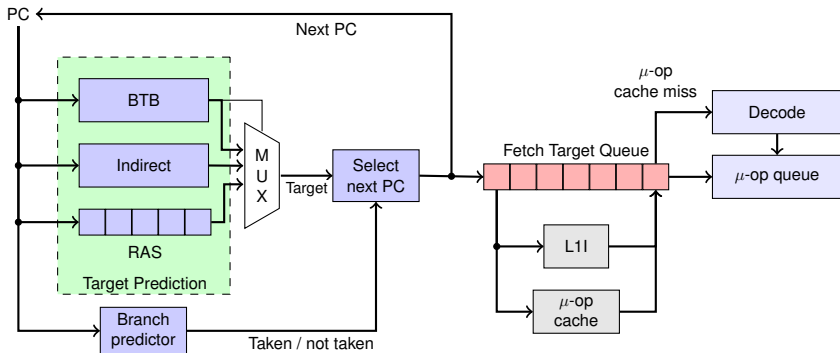
# BACKGROUND & MOTIVATION

## PROCESSOR FRONT-END



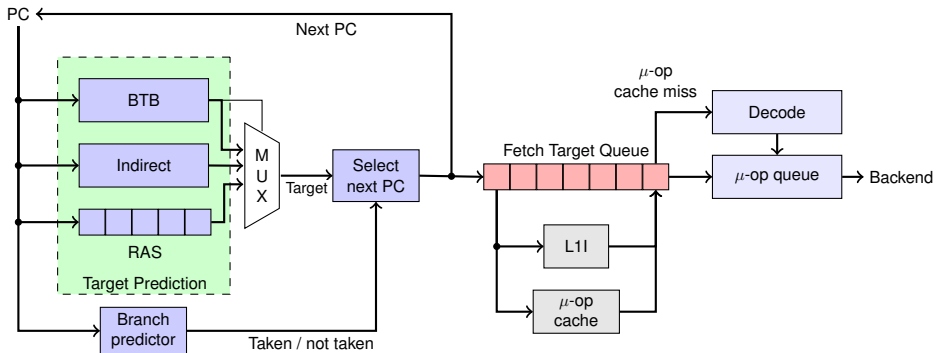
# BACKGROUND & MOTIVATION

## PROCESSOR FRONT-END



# BACKGROUND & MOTIVATION

## PROCESSOR FRONT-END

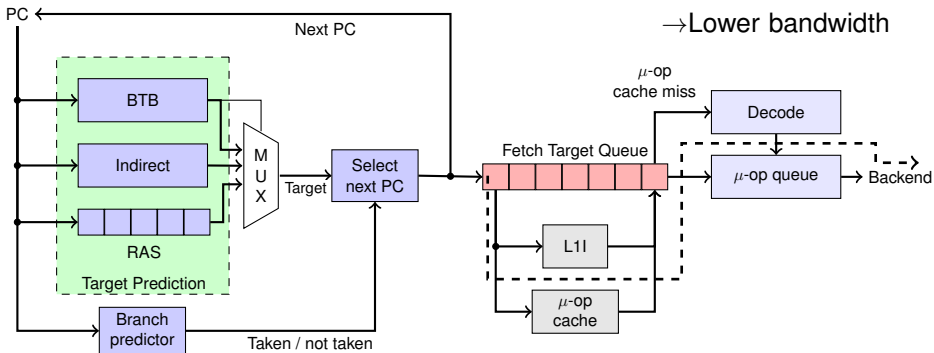


# BACKGROUND & MOTIVATION

## PROCESSOR FRONT-END

- Decode latency
- Decode energy
- Lower bandwidth

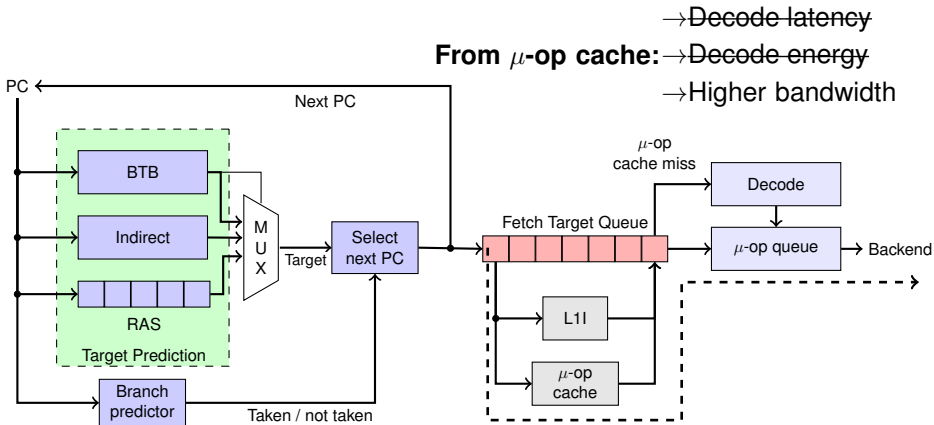
**From L1I:**





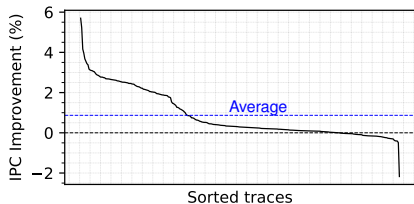
# BACKGROUND & MOTIVATION

## PROCESSOR FRONT-END



# BACKGROUND & MOTIVATION

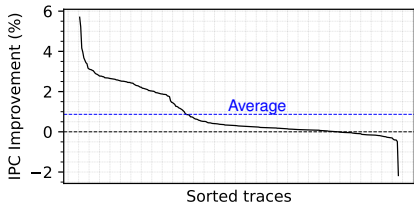
## PERFORMANCE OF $\mu$ -OPS CACHE WITH SERVER WORKLOADS



# BACKGROUND & MOTIVATION

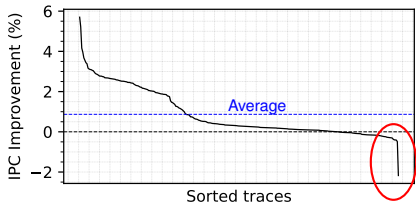
## PERFORMANCE OF $\mu$ -OPS CACHE WITH SERVER WORKLOADS

→ 4Kops  $\mu$ -op provides only **0.87%** improvement over no  $\mu$ -op cache



# BACKGROUND & MOTIVATION

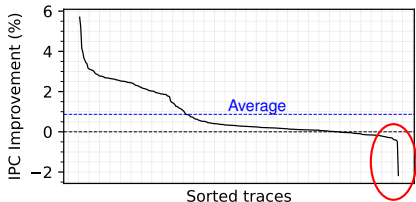
## PERFORMANCE OF $\mu$ -OPS CACHE WITH SERVER WORKLOADS



- 4Kops  $\mu$ -op provides only **0.87%** improvement over no  $\mu$ -op cache
- **19.3%** of traces show a slowdown

# BACKGROUND & MOTIVATION

## PERFORMANCE OF $\mu$ -OPS CACHE WITH SERVER WORKLOADS

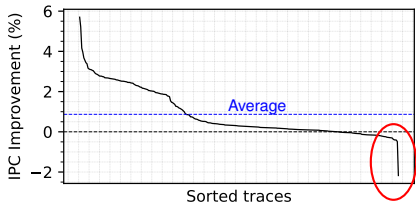


- 4Kops  $\mu$ -op provides only **0.87%** improvement over no  $\mu$ -op cache
- **19.3%** of traces show a slowdown
- Increasing **size** of  $\mu$ -op cache does not help

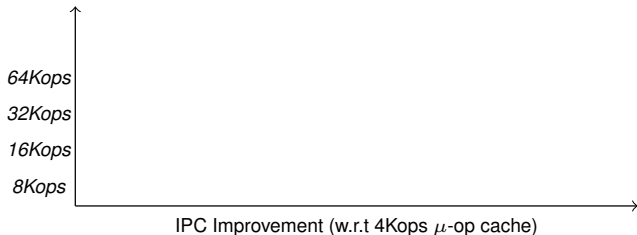
→  
IPC Improvement (w.r.t 4Kops  $\mu$ -op cache)

# BACKGROUND & MOTIVATION

## PERFORMANCE OF $\mu$ -OPS CACHE WITH SERVER WORKLOADS

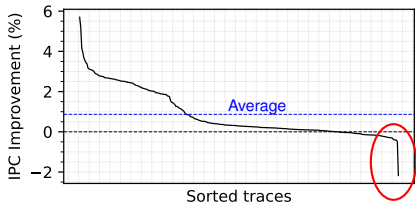


- 4Kops  $\mu$ -op provides only **0.87%** improvement over no  $\mu$ -op cache
- **19.3%** of traces show a slowdown
- Increasing **size** of  $\mu$ -op cache does not help

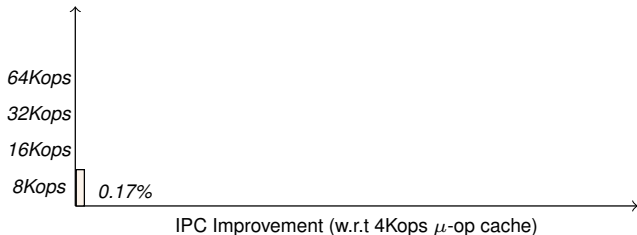


# BACKGROUND & MOTIVATION

## PERFORMANCE OF $\mu$ -OPS CACHE WITH SERVER WORKLOADS

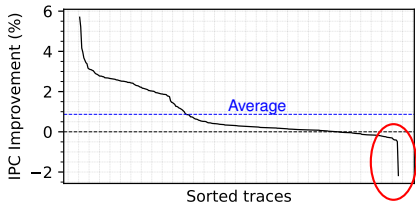


- 4Kops  $\mu$ -op provides only **0.87%** improvement over no  $\mu$ -op cache
- **19.3%** of traces show a slowdown
- Increasing **size** of  $\mu$ -op cache does not help

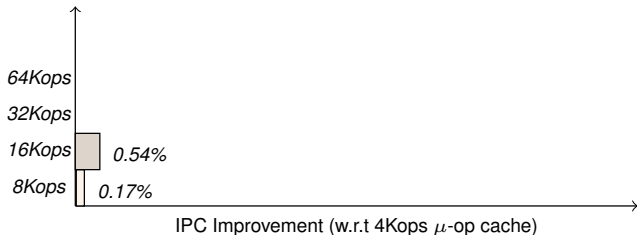


# BACKGROUND & MOTIVATION

## PERFORMANCE OF $\mu$ -OPS CACHE WITH SERVER WORKLOADS



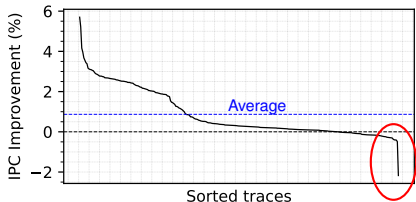
- 4Kops  $\mu$ -op provides only **0.87%** improvement over no  $\mu$ -op cache
- **19.3%** of traces show a slowdown
- Increasing **size** of  $\mu$ -op cache does not help



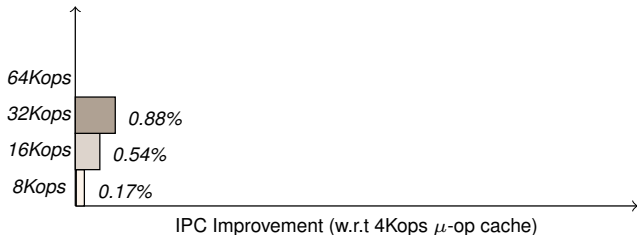


# BACKGROUND & MOTIVATION

## PERFORMANCE OF $\mu$ -OPS CACHE WITH SERVER WORKLOADS

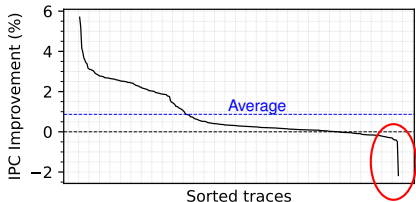


- 4Kops  $\mu$ -op provides only **0.87%** improvement over no  $\mu$ -op cache
- **19.3%** of traces show a slowdown
- Increasing **size** of  $\mu$ -op cache does not help

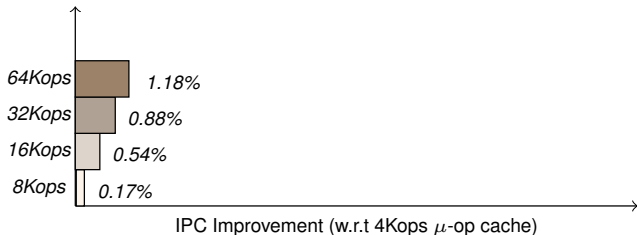


# BACKGROUND & MOTIVATION

## PERFORMANCE OF $\mu$ -OPS CACHE WITH SERVER WORKLOADS

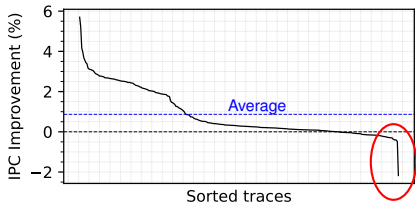


- 4Kops  $\mu$ -op provides only **0.87%** improvement over no  $\mu$ -op cache
- **19.3%** of traces show a slowdown
- Increasing **size** of  $\mu$ -op cache does not help

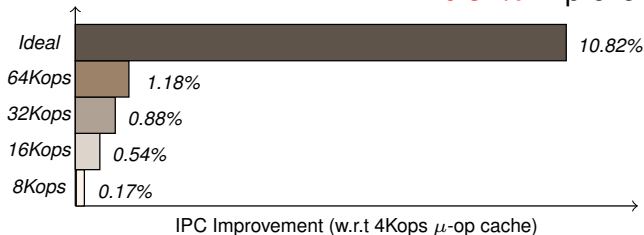


# BACKGROUND & MOTIVATION

## PERFORMANCE OF $\mu$ -OPS CACHE WITH SERVER WORKLOADS

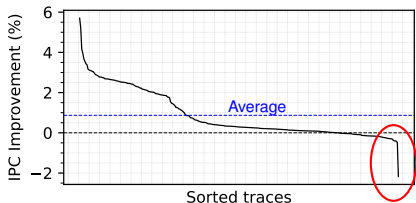


- 4Kops  $\mu$ -op provides only **0.87%** improvement over no  $\mu$ -op cache
- **19.3%** of traces show a slowdown
- Increasing **size** of  $\mu$ -op cache does not help
- Ideal  $\mu$ -op cache can provide **10.82%** improvement

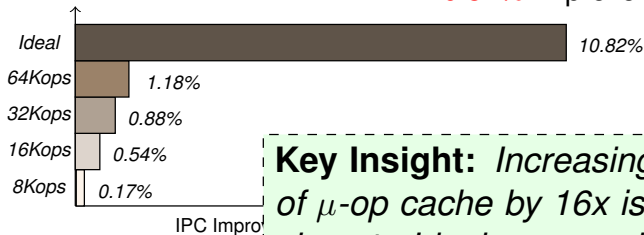


# BACKGROUND & MOTIVATION

## PERFORMANCE OF $\mu$ -OPS CACHE WITH SERVER WORKLOADS



- 4Kops  $\mu$ -op provides only **0.87%** improvement over no  $\mu$ -op cache
- **19.3%** of traces show a slowdown
- Increasing **size** of  $\mu$ -op cache does not help
- Ideal  $\mu$ -op cache can provide **10.82%** improvement



**Key Insight:** *Increasing the size of  $\mu$ -op cache by 16x is still not close to Ideal  $\mu$ -op cache*

## BACKGROUND & MOTIVATION

### FTQ BEHAVIOR ON BRANCH MISS

- FTQ is unable to hide the L1I miss latency on **branch misprediction**

# BACKGROUND & MOTIVATION

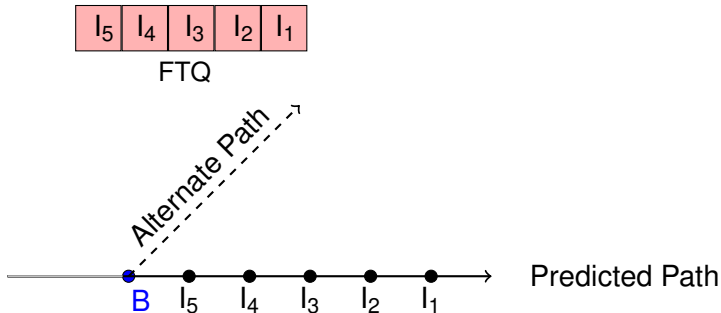
## FTQ BEHAVIOR ON BRANCH MISS

- FTQ is unable to hide the L1I miss latency on **branch misprediction**
  - **FTQ is cleared** on a branch misprediction

# BACKGROUND & MOTIVATION

## FTQ BEHAVIOR ON BRANCH MISS

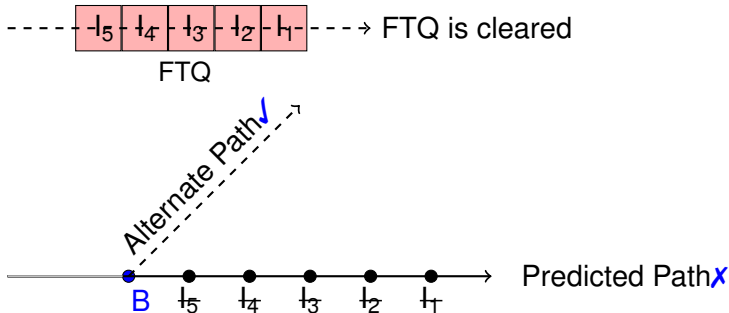
- FTQ is unable to hide the L1I miss latency on **branch misprediction**
- **FTQ is cleared** on a branch misprediction



# BACKGROUND & MOTIVATION

## FTQ BEHAVIOR ON BRANCH MISS

- FTQ is unable to hide the L1I miss latency on **branch misprediction**
  - **FTQ is cleared** on a branch misprediction





## BACKGROUND & MOTIVATION

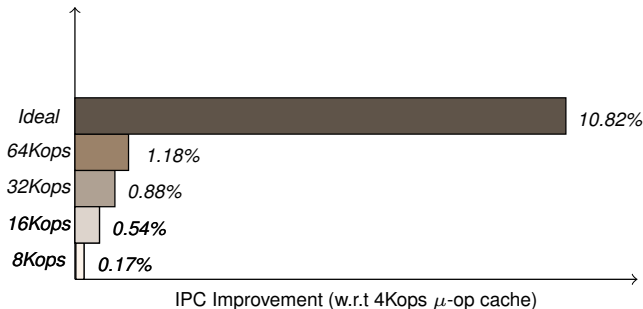
### FTQ BEHAVIOR ON BRANCH MISS

- FTQ is unable to hide the L1I miss latency on **branch misprediction**
  - **FTQ is cleared** on a branch misprediction
- What if the **correct path** was always **in the  $\mu$ -op cache** after a pipeline flush due to branch misprediction?

## BACKGROUND & MOTIVATION

### FTQ BEHAVIOR ON BRANCH MISS

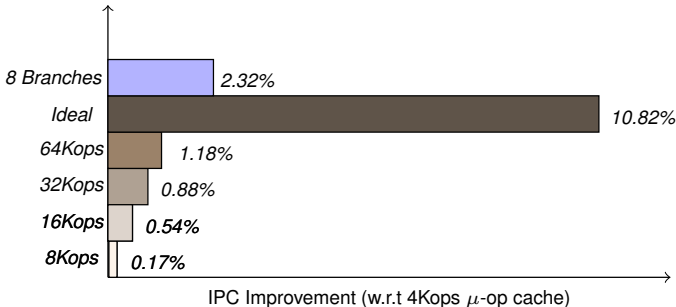
- FTQ is unable to hide the L1I miss latency on **branch misprediction**
  - **FTQ is cleared** on a branch misprediction
- What if the **correct path** was always **in the  $\mu$ -op cache** after a pipeline flush due to branch misprediction?



## BACKGROUND & MOTIVATION

### FTQ BEHAVIOR ON BRANCH MISS

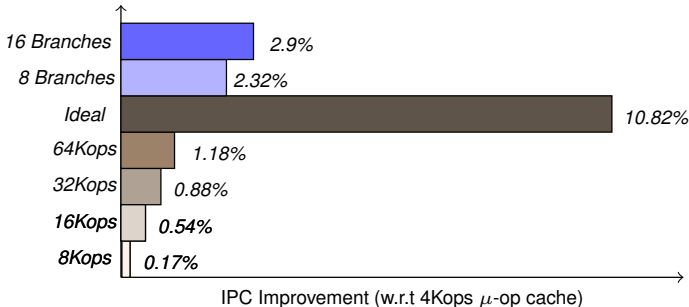
- FTQ is unable to hide the L1I miss latency on **branch misprediction**
  - **FTQ is cleared** on a branch misprediction
- What if the **correct path** was always **in the  $\mu$ -op cache** after a pipeline flush due to branch misprediction?



# BACKGROUND & MOTIVATION

## FTQ BEHAVIOR ON BRANCH MISS

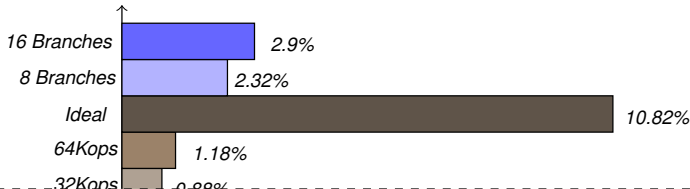
- FTQ is unable to hide the L1I miss latency on **branch misprediction**
  - **FTQ is cleared** on a branch misprediction
- What if the **correct path** was always **in the  $\mu$ -op cache** after a pipeline flush due to branch misprediction?



# BACKGROUND & MOTIVATION

## FTQ BEHAVIOR ON BRANCH MISS

- FTQ is unable to hide the L1I miss latency on **branch misprediction**
  - **FTQ is cleared** on a branch misprediction
- What if the **correct path** was always **in the  $\mu$ -op cache** after a pipeline flush due to branch misprediction?

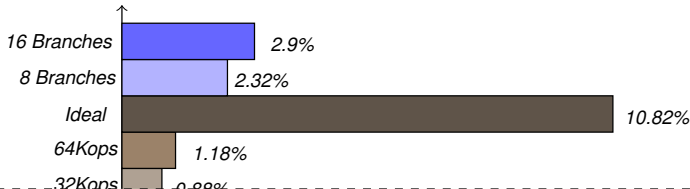


**Key Insight:**

# BACKGROUND & MOTIVATION

## FTQ BEHAVIOR ON BRANCH MISS

- FTQ is unable to hide the L1I miss latency on **branch misprediction**
  - **FTQ is cleared** on a branch misprediction
- What if the **correct path** was always **in the  $\mu$ -op cache** after a pipeline flush due to branch misprediction?



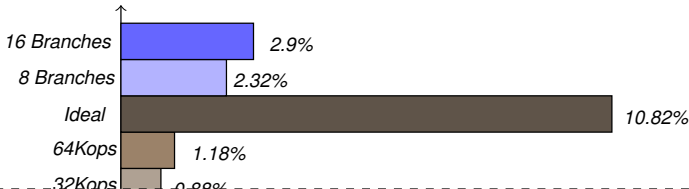
### Key Insight:

1. *FTQ is unable to hide the fetch latency on branch misprediction*

## BACKGROUND & MOTIVATION

### FTQ BEHAVIOR ON BRANCH MISS

- FTQ is unable to hide the L1I miss latency on **branch misprediction**
  - **FTQ is cleared** on a branch misprediction
- What if the **correct path** was always **in the  $\mu$ -op cache** after a pipeline flush due to branch misprediction?



#### Key Insight:

1. FTQ is unable to hide the fetch latency on branch misprediction
2. Focusing on few but critical instructions can provide significant performance benefits

# OUTLINE

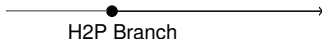
- 1 Overview
- 2 Background & Motivation
- 3 UCP
- 4 Methodology & Results
- 5 Conclusions



# UCP

## UCP: OVERVIEW

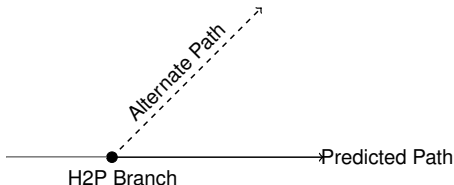
- ① Identifies a **hard-to-predict** conditional branch (H2P)



# UCP

## UCP: OVERVIEW

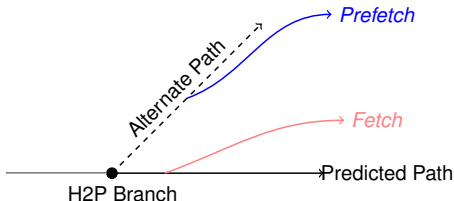
- ① Identifies a **hard-to-predict** conditional branch (H2P)
- ② On a H2P begin generating **addresses on alternate path** (alternate path)



# UCP

## UCP: OVERVIEW

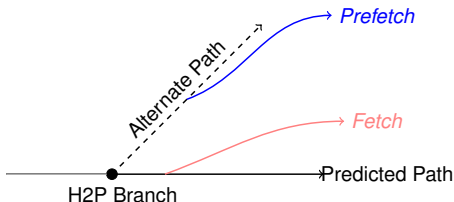
- ① Identifies a **hard-to-predict** conditional branch (H2P)
- ② On a H2P begin generating **addresses on alternate path** (alternate path)
- ③ Generated addresses on alternate path are **prefetched** in parallel to predicted path fetch



# UCP

## UCP: OVERVIEW

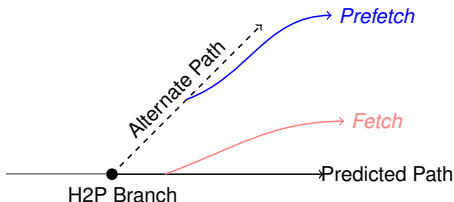
- ① Identifies a **hard-to-predict** conditional branch (H2P)
- ② On a H2P begin generating **addresses on alternate path** (alternate path)
- ③ Generated addresses on alternate path are **prefetched** in parallel to predicted path fetch
- ④ Prefetched instructions are decoded and **inserted in the  $\mu$ -op cache**



# UCP

## UCP: OVERVIEW

- ① Identifies a **hard-to-predict** conditional branch (H2P)
- ② On a H2P begin generating **addresses on alternate path** (alternate path)
- ③ Generated addresses on alternate path are **prefetched** in parallel to predicted path fetch
- ④ Prefetched instructions are decoded and **inserted in the  $\mu$ -op cache**



**Key Idea:** *Keep the alternate path in the  $\mu$ -op cache for H2P branches*

# UCP

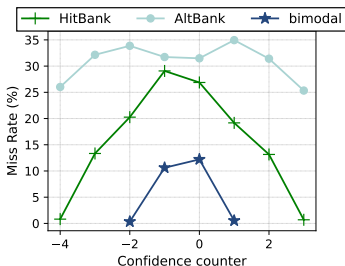
## UCP: H2P BRANCH DETECTION

→ **H2P Branch:** a branch which has high chance of being mispredicted

# UCP

## UCP: H2P BRANCH DETECTION

- **H2P Branch:** a branch which has high chance of being mispredicted
  - TAGE-Conf<sup>2</sup>



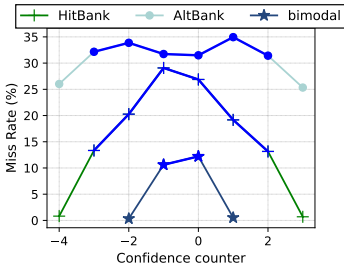
# UCP

## UCP: H2P BRANCH DETECTION

→ **H2P Branch**: a branch which has high chance of being mispredicted

→ TAGE-Conf<sup>2</sup>

- **Not saturated predictions** from AltBank, HitBank & BiModal





# UCP

## UCP: H2P BRANCH DETECTION

- **H2P Branch**: a branch which has high chance of being mispredicted
  - TAGE-Conf<sup>2</sup>
    - **Not saturated predictions** from AltBank, HitBank & BiModal
    - Does not consider **SC and LP**

# UCP

## UCP: H2P BRANCH DETECTION

- **H2P Branch:** a branch which has high chance of being mispredicted
  - TAGE-Conf<sup>2</sup>
    - **Not saturated predictions** from AltBank, HitBank & BiModal
    - Does not consider **SC and LP**
  - UCP-Conf

# UCP

## UCP: H2P BRANCH DETECTION

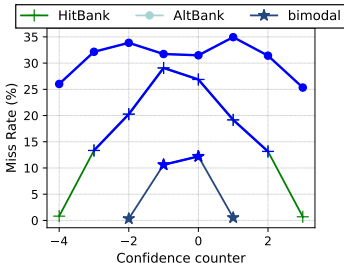
→ **H2P Branch**: a branch which has high chance of being mispredicted

→ TAGE-Conf<sup>2</sup>

- **Not saturated predictions** from AltBank, HitBank & BiModal
- Does not consider **SC and LP**

→ UCP-Conf

- All predictions from **AltBanks** shows high miss rate



# UCP

## UCP: H2P BRANCH DETECTION

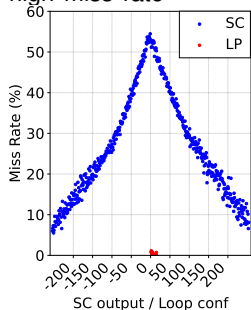
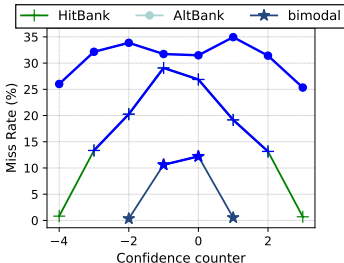
→ **H2P Branch**: a branch which has high chance of being mispredicted

→ TAGE-Conf<sup>2</sup>

- **Not saturated predictions** from AltBank, HitBank & BiModal
- Does not consider **SC and LP**

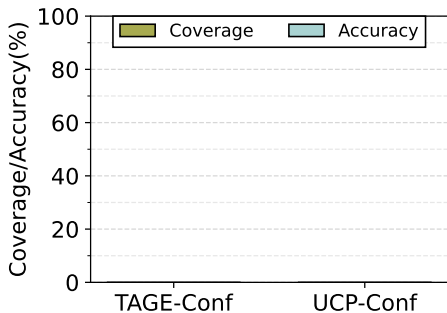
→ UCP-Conf

- All predictions from **AltBanks** shows high miss rate
- **SC** shows high miss rate



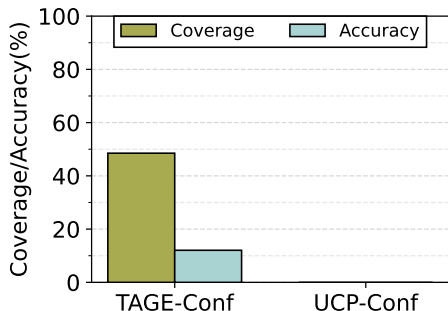
# UCP

## UCP: H2P BRANCH DETECTION



# UCP

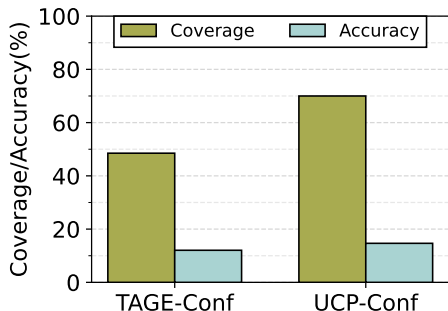
## UCP: H2P BRANCH DETECTION



→ TAGE-Conf provide 48.5% coverage and 12% accuracy

# UCP

## UCP: H2P BRANCH DETECTION

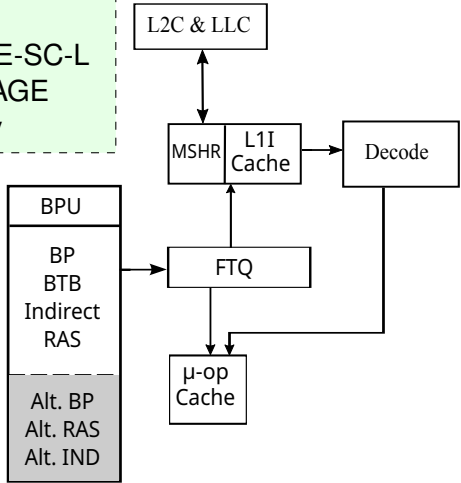


- TAGE-Conf provide 48.5% coverage and 12% accuracy
- UCP-Conf **improve coverage** to 70% **and accuracy** to 14.66%

# UCP

## UCP: DESIGN

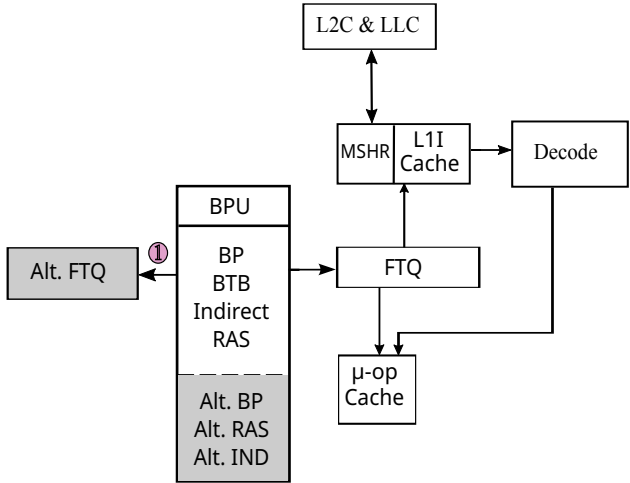
**BTB:** Shared  
**Alt BP:** 8KB TAGE-SC-L  
**Alt IND:** 4KB ITTAGE  
**Alt RAS:** 16-entry





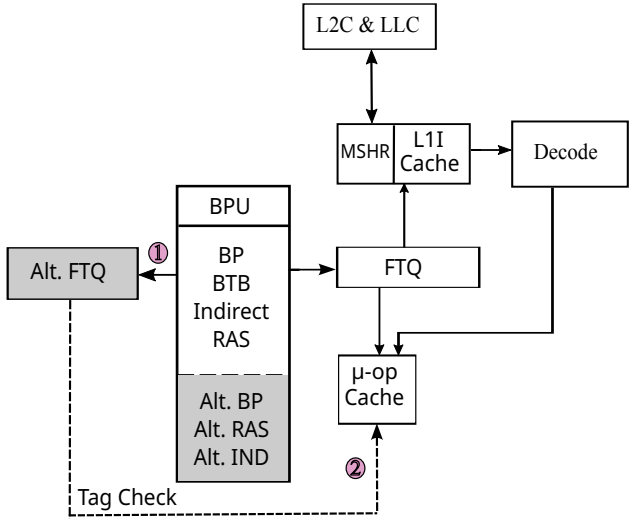
# UCP

## UCP: DESIGN



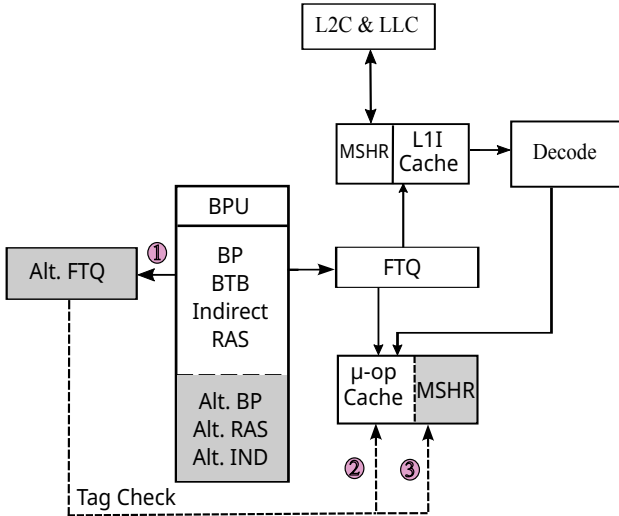
# UCP

## UCP: DESIGN



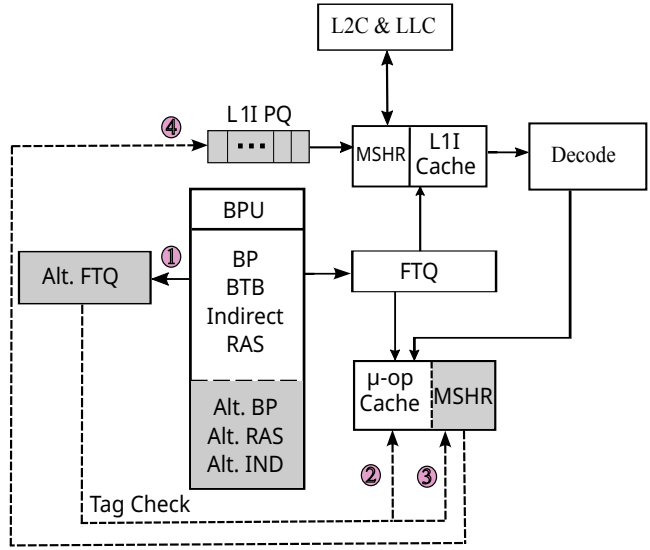
# UCP

## UCP: DESIGN



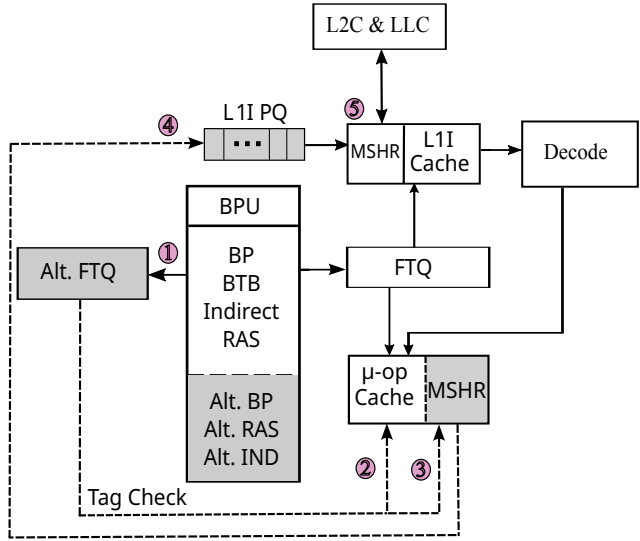
# UCP

## UCP: DESIGN



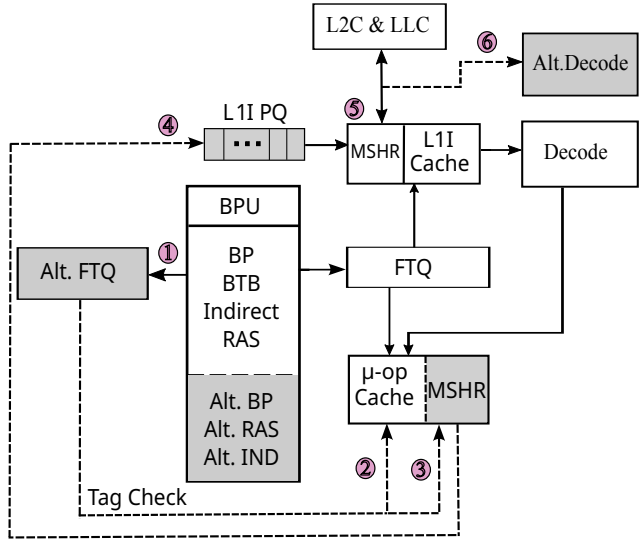
# UCP

## UCP: DESIGN



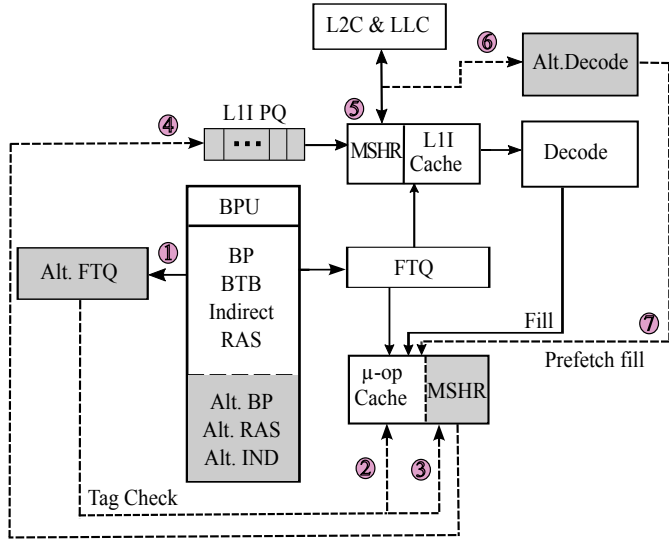
# UCP

## UCP: DESIGN



# UCP

## UCP: DESIGN



# UCP

## UCP: ALTERNATE PATH STOPPING CONDITIONS

→ Stopping alternate path is crucial to prevent  $\mu$ -op cache pollution



# UCP

## UCP: ALTERNATE PATH STOPPING CONDITIONS

- Stopping alternate path is crucial to prevent  $\mu$ -op cache pollution
- When a new H2P branch is detected

# UCP

## UCP: ALTERNATE PATH STOPPING CONDITIONS

- Stopping alternate path is crucial to prevent  $\mu$ -op cache pollution
- When a new H2P branch is detected
- When the alternate path is unlikely to become the correct path

## UCP

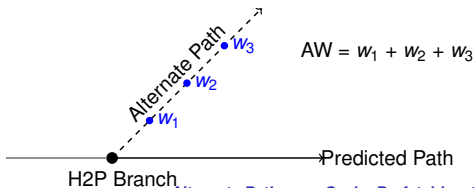
### UCP: ALTERNATE PATH STOPPING CONDITIONS

- Stopping alternate path is crucial to prevent  $\mu$ -op cache pollution
- When a new H2P branch is detected
- When the alternate path is unlikely to become the correct path
  - BTB miss on predicted taken branch on the alternate path

# UCP

## UCP: ALTERNATE PATH STOPPING CONDITIONS

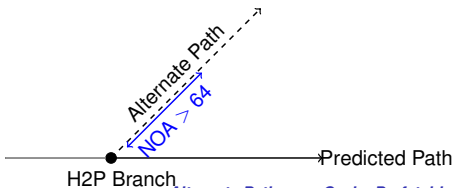
- Stopping alternate path is crucial to prevent  $\mu$ -op cache pollution
- When a **new H2P** branch is detected
- When the alternate path is **unlikely** to become the correct path
  - **BTB miss** on predicted taken branch on the alternate path
  - Weight of **each branch** on the alternate path is **accumulated**, high confidence branches have lower weight.



# UCP

## UCP: ALTERNATE PATH STOPPING CONDITIONS

- Stopping alternate path is crucial to prevent  **$\mu$ -op cache pollution**
- When a **new H2P** branch is detected
- When the alternate path is **unlikely** to become the correct path
  - **BTB miss** on predicted taken branch on the alternate path
  - Weight of **each branch** on the alternate path is **accumulated**, high confidence branches have lower weight.
  - **Non-branch instructions** after a branch are counted. Once the count reaches 64 alternate paths **stops** in our work



# OUTLINE

- 1 Overview
- 2 Background & Motivation
- 3 UCP
- 4 **Methodology & Results**
- 5 Conclusions

# METHODOLOGY & RESULTS

## SIMULATION SETUP

- **ChampSim** + subset (traces showing  $\geq 5\%$  improvement with ideal  $\mu$ -op cache) of **CVP traces** (2 FP, 97 INT, 73 Crypto and 134 datacenter trace)

# METHODOLOGY & RESULTS

## SIMULATION SETUP

- **ChampSim** + subset (traces showing  $\geq 5\%$  improvement with ideal  $\mu$ -op cache) of **CVP traces** (2 FP, 97 INT, 73 Crypto and 134 datacenter trace)
- Intel **Alder Lake** like microarchitecture



# METHODOLOGY & RESULTS

## SIMULATION SETUP

- **ChampSim** + subset (traces showing  $\geq 5\%$  improvement with ideal  $\mu$ -op cache) of **CVP traces** (2 FP, 97 INT, 73 Crypto and 134 datacenter trace)
- Intel **Alder Lake** like microarchitecture
- We execute 100M instructions, **50M warmup** and **50M to collect stats**

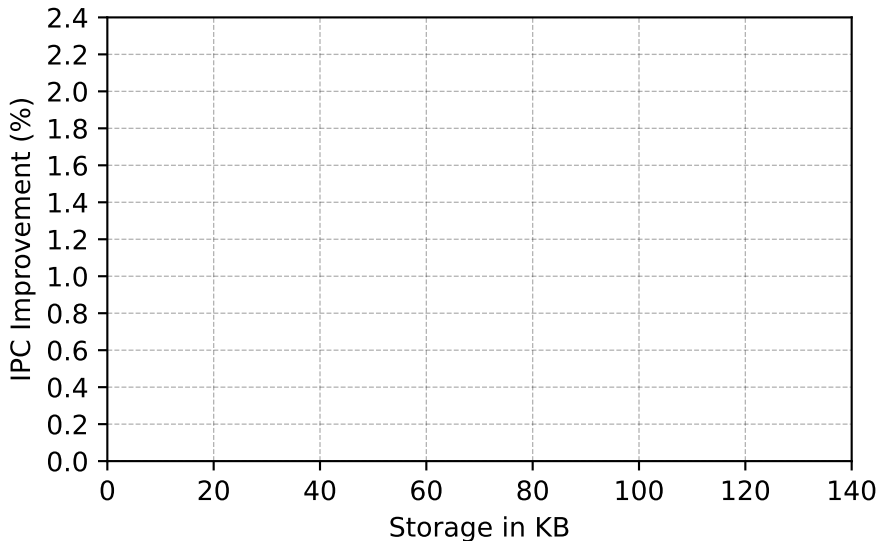
# METHODOLOGY & RESULTS

## SIMULATION SETUP

- **ChampSim** + subset (traces showing  $\geq 5\%$  improvement with ideal  $\mu$ -op cache) of **CVP traces** (2 FP, 97 INT, 73 Crypto and 134 datacenter trace)
- Intel **Alder Lake** like microarchitecture
- We execute 100M instructions, **50M warmup** and **50M to collect stats**
- **1 cycle penalty** for switching from  $\mu$ -op cache to L1I cache

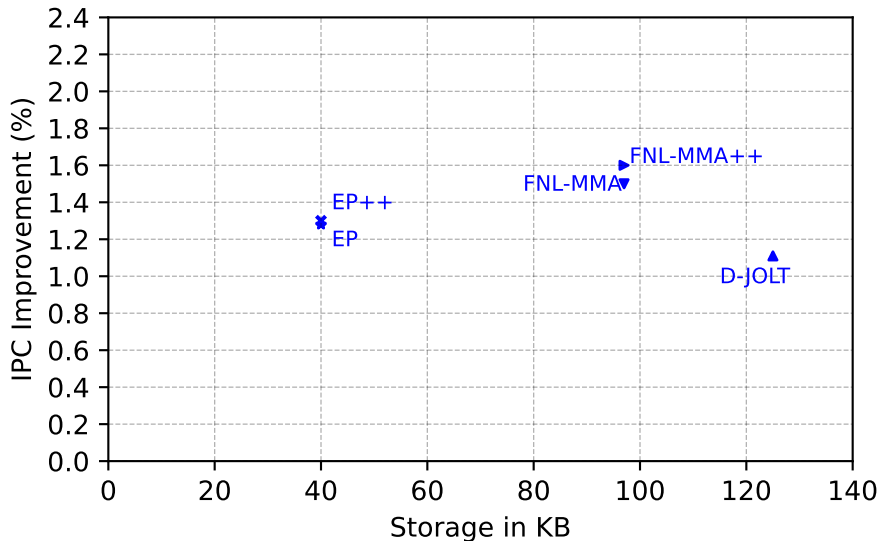
# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



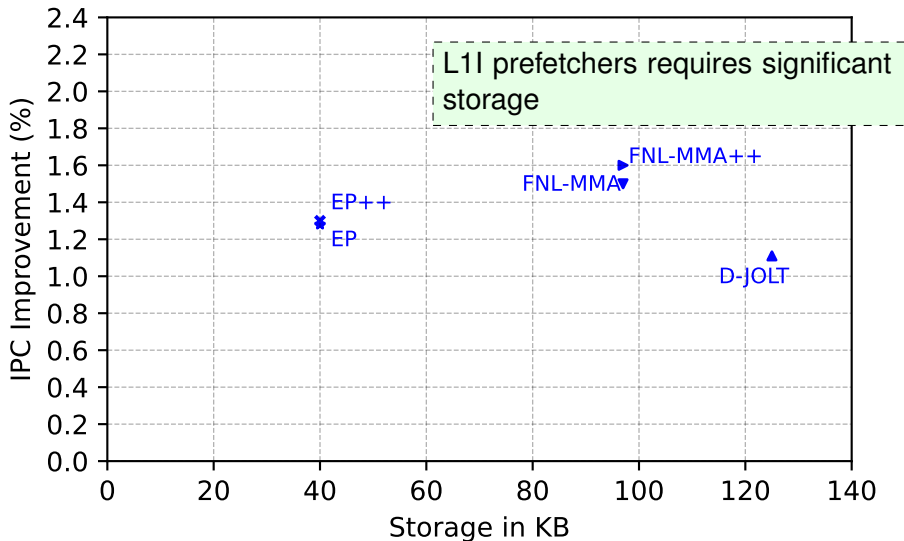
# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



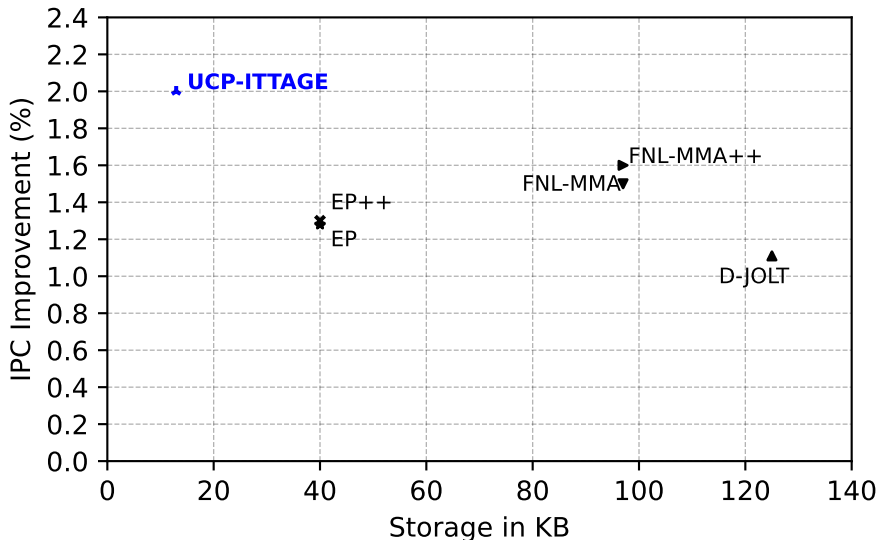
# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



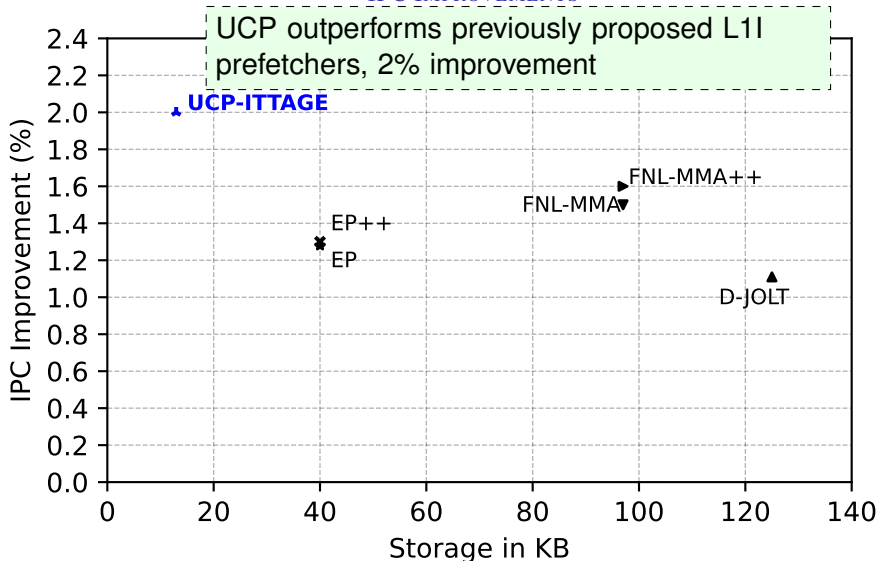
# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



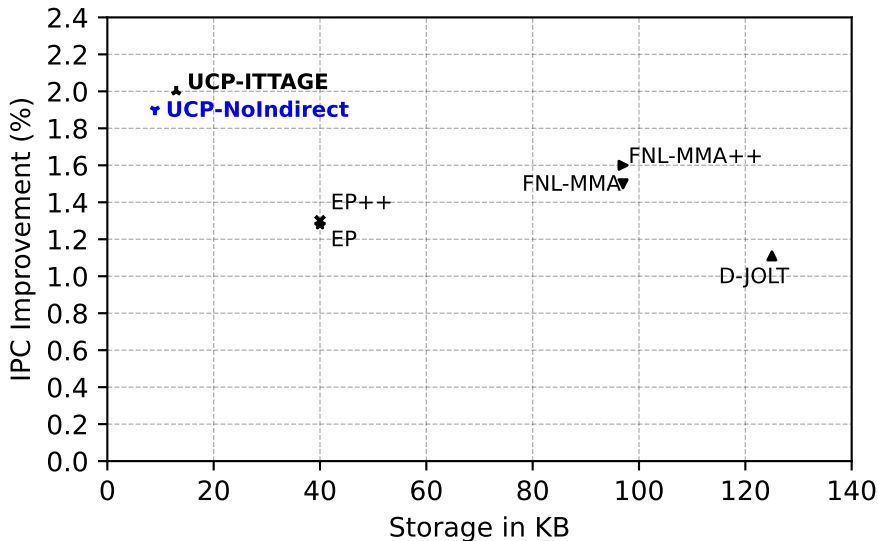
# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



# METHODOLOGY & RESULTS

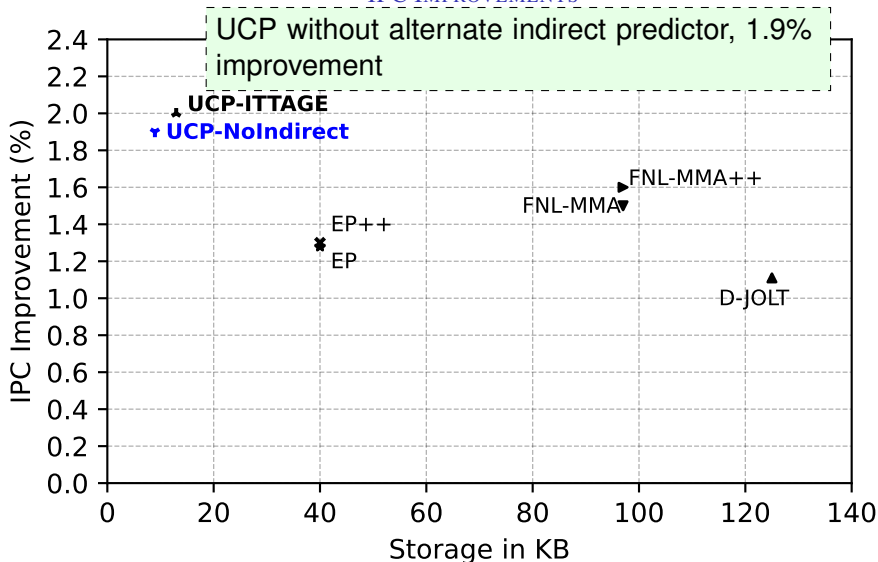
## IPC IMPROVEMENTS





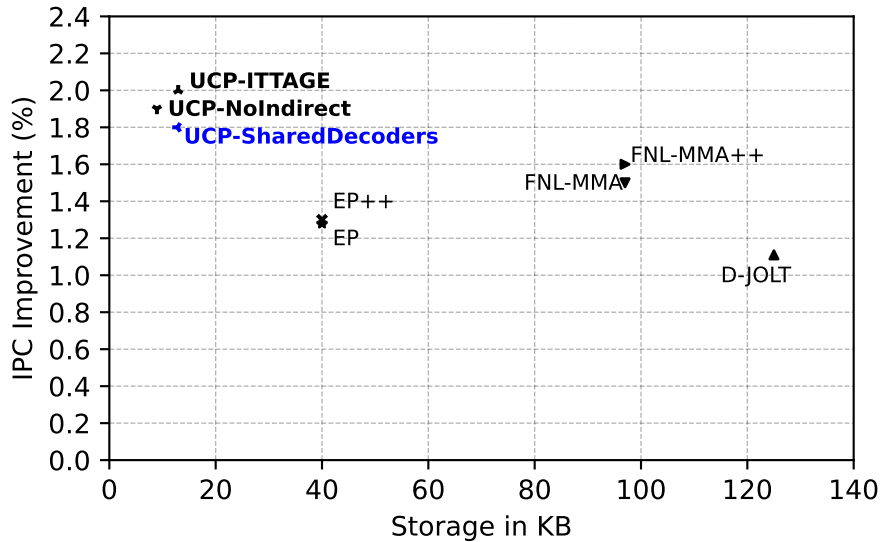
# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



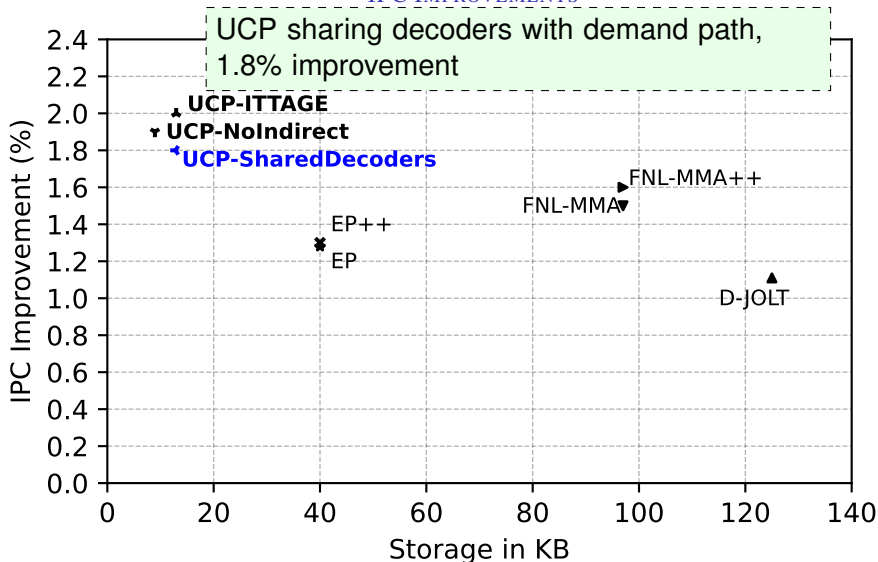
# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



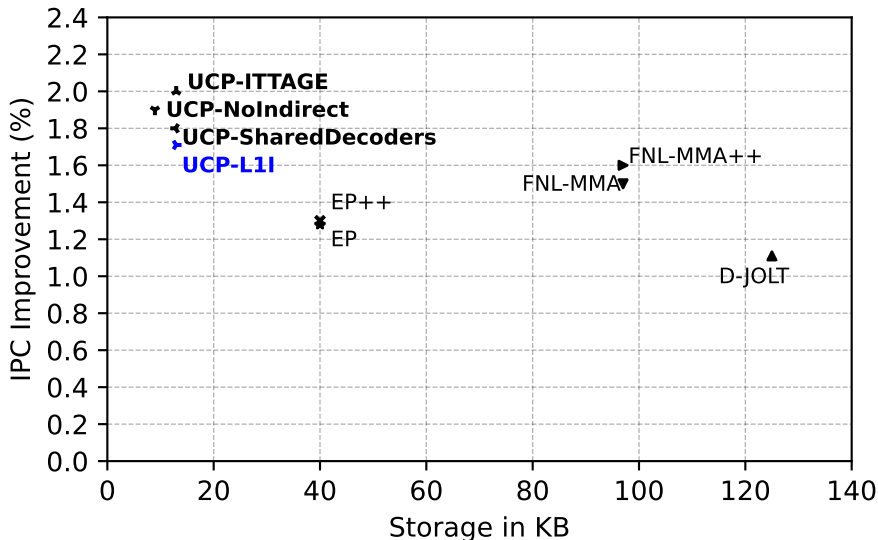
# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



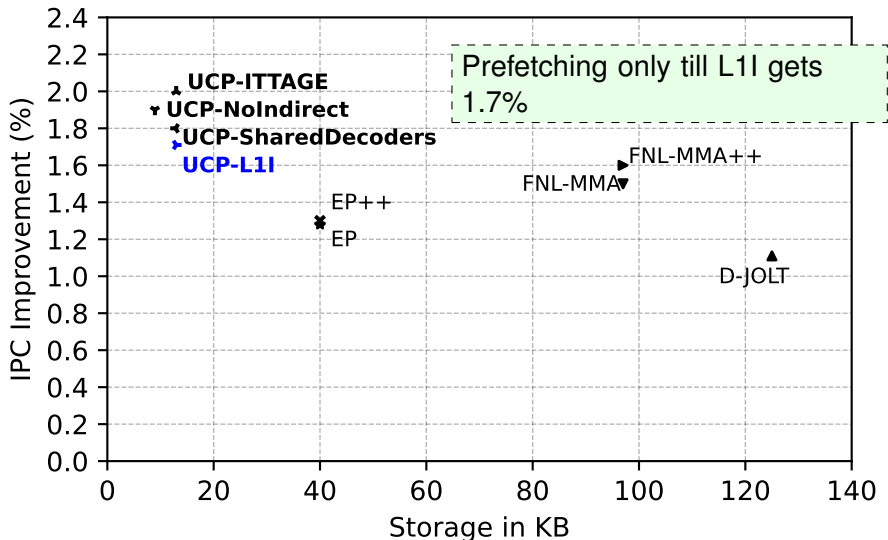
# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



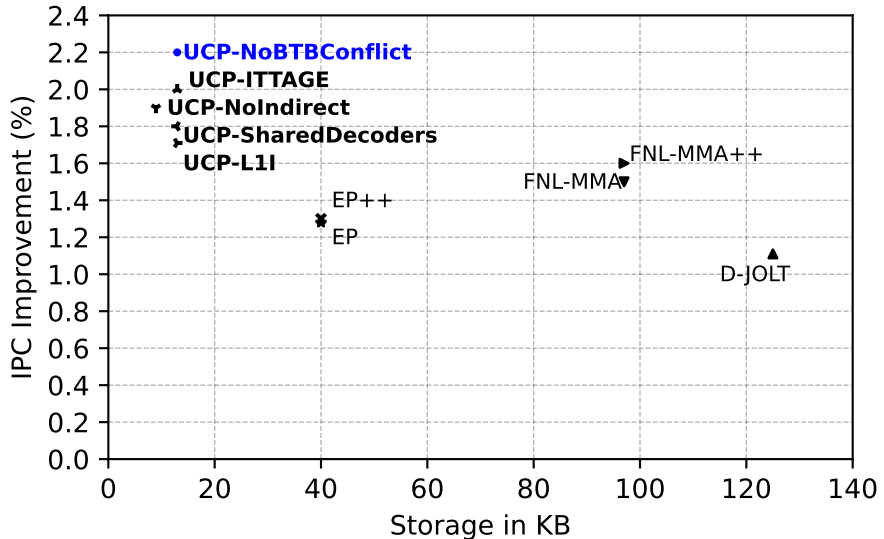
# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



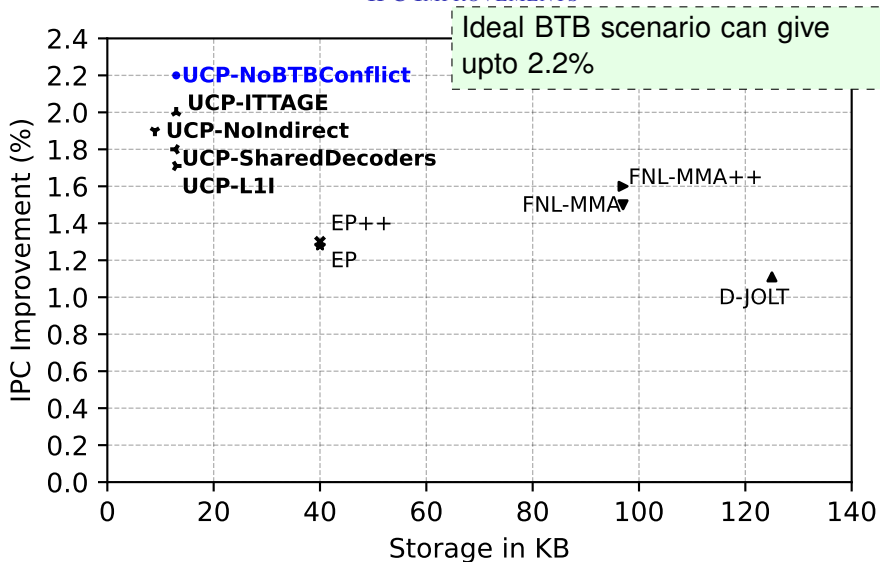
# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



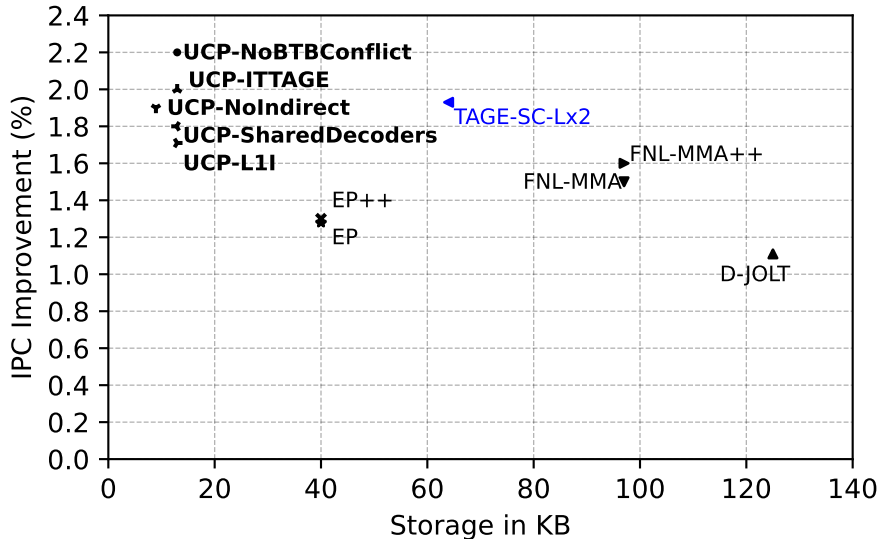
# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



# METHODOLOGY & RESULTS

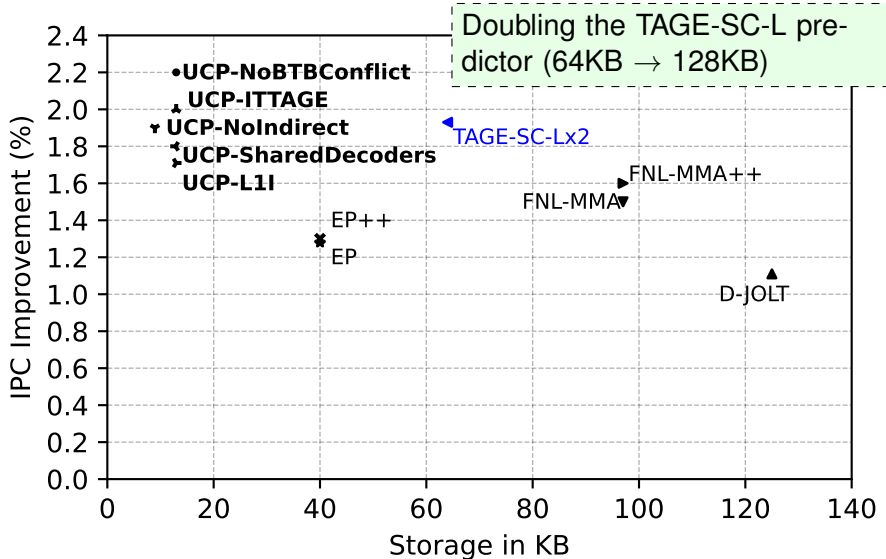
## IPC IMPROVEMENTS





# METHODOLOGY & RESULTS

## IPC IMPROVEMENTS



# OUTLINE

- 1 Overview
- 2 Background & Motivation
- 3 UCP
- 4 Methodology & Results
- 5 **Conclusions**

## CONCLUSIONS

→ FTQ fails to hide **L1 miss latency** on branch miss

## CONCLUSIONS

- FTQ fails to hide **L1 miss latency** on branch miss
- Focusing only **a few but critical** instructions can provide better performance

## CONCLUSIONS

- FTQ fails to hide **L1 miss latency** on branch miss
- Focusing only **a few but critical** instructions can provide better performance
- UCP focus on critical instructions after a **H2P branch**

## CONCLUSIONS

- FTQ fails to hide **L1I miss latency** on branch miss
- Focusing only **a few but critical** instructions can provide better performance
- UCP focus on critical instructions after a **H2P branch**
- Still space for improvement in **optimizing  $\mu$ -op cache**

# ALTERNATE PATH $\mu$ -OP CACHE PREFETCHING

Sawan Singh<sup>1</sup> Arthur Perais<sup>2</sup> Alexandra Jimborean<sup>1</sup>  
Alberto Ros<sup>1</sup>



singh.sawan@um.es

Thank you for your attention!



ECHO, ERC Consolidator Grant (No 819134)

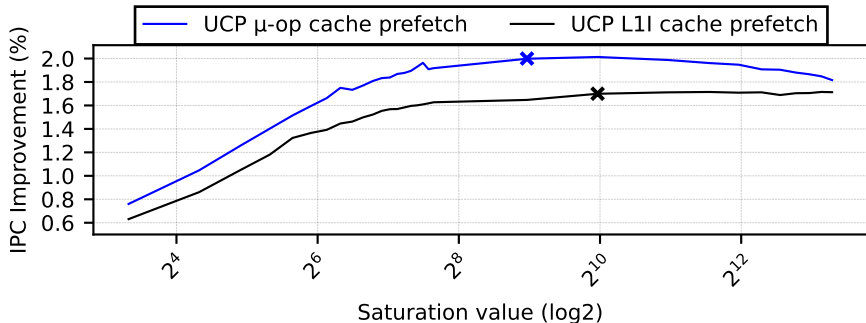
This presentation and recording belong to the authors. No distribution is allowed without the authors' permission.

## BACKUP SLIDES[BTB BANKING]

- UCP reuses the BTB by **doubling the number of BTB banks** (from 16 to 32)
- Each cycle we determine the **banks to be accessed**
- By default, **demand requests are given priority** to access the conflict banks
- UCP keeps a 3-bit saturated counter which is incremented every time the **alternate path is delayed**
- When the **counter saturates**, the alternate path is given priority for the conflict banks in that cycle
- The counter **resets next cycle**



# BACKUP SLIDES[SATURATION COUNTER]



## BACKUP SLIDES[L1I PREFETCHERS]

