# ALTERNATE PATH µ-OP CACHE PREFETCHING

**Sawan Singh**[1]    Arthur Perais[2]    Alexandra Jimborean[1]
Alberto Ros[1]

[1]Computer Engineering Department
University of Murcia, Spain

[2]TIMA, Univ. Grenoble Alpes, CNRS,
Grenoble INP, Grenoble, France

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

$\rightarrow$ $\mu$-op Cache

$\rightarrow$ $\mu$-op Cache
- Holds recently decoded $\mu$-ops

# OVERVIEW

$\rightarrow$ $\mu$-op Cache
- Holds recently decoded $\mu$-ops
- First introduced for energy savings[1] in x86 which requires complex decoders

[1] *Solomon et al. Micro-operation cache: a power aware frontend for variable instruction length ISA*

→ $\mu$-op Cache

- Holds recently decoded $\mu$-ops
- First introduced for energy savings[1] in x86 which requires complex decoders
- Potential to provide an IPC improvement of 10.82%

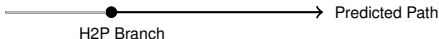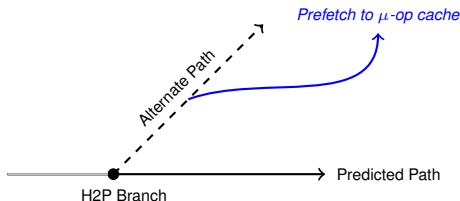[1] *Solomon et al. Micro-operation cache: a power aware frontend for variable instruction length ISA*

→ $\mu$-op Cache
- Holds recently decoded $\mu$-ops
- First introduced for energy savings[1] in x86 which requires complex decoders
- Potential to provide an IPC improvement of 10.82%

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

[1] *Solomon et al. Micro-operation cache: a power aware frontend for variable instruction length ISA*

$\rightarrow$ $\mu$-op Cache
  - Holds recently decoded $\mu$-ops
  - First introduced for energy savings[1] in x86 which requires complex decoders
  - Potential to provide an IPC improvement of 10.82%

$\rightarrow$ We propose UCP (*Alternate Path $\mu$-op Cache Prefetching*)
  - ① Identify hard-to-predict branches



H2P Branch → Predicted Path

[1] *Solomon et al. Micro-operation cache: a power aware frontend for variable instruction length ISA*

→ $\mu$-op Cache
  - Holds recently decoded $\mu$-ops
  - First introduced for energy savings[1] in x86 which requires complex decoders
  - Potential to provide an IPC improvement of 10.82%

→ We propose UCP (*Alternate Path $\mu$-op Cache Prefetching*)
  ① Identify hard-to-predict branches
  ② Prefetch $\mu$-ops from alternate path



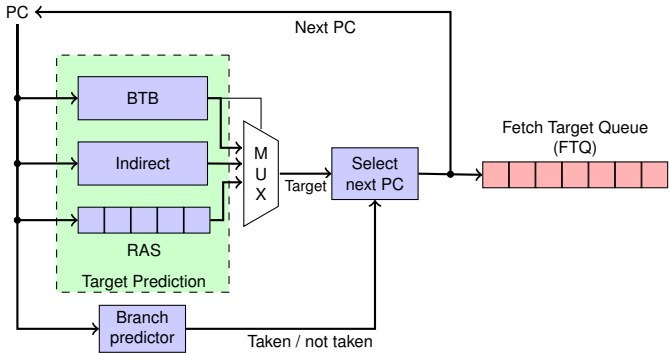*Prefetch to $\mu$-op cache*

Alternate Path

Predicted Path

H2P Branch

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

From L1I:
→Decode latency
→Decode energy
→Lower bandwidth

*Alternate Path μ-op Cache Prefetching @ISCA'51*

From $\mu$-op cache:
→Decode latency
→Decode energy
→Higher bandwidth

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

$\rightarrow$ Server workloads overwhelm current $\mu$-op caches

$\quad \rightarrow 0.87\%$ IPC improvement over no $\mu$-op cache

$\rightarrow$ Server workloads overwhelm current $\mu$-op caches

    $\rightarrow$ $0.87\%$ IPC improvement over no $\mu$-op cache

$\rightarrow$ Increasing size of $\mu$-op cache does not help

$\rightarrow$ Server workloads overwhelm current $\mu$-op caches
  $\rightarrow$ $0.87\%$ IPC improvement over no $\mu$-op cache
$\rightarrow$ Increasing size of $\mu$-op cache does not help

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

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

$\rightarrow$ Server workloads overwhelm current $\mu$-op caches
  $\rightarrow$ $0.87\%$ IPC improvement over no $\mu$-op cache
$\rightarrow$ Increasing size of $\mu$-op cache does not help

*64Kops*

*32Kops*

*16Kops*

*8Kops*

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

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

$\rightarrow$ Server workloads overwhelm current $\mu$-op caches
   $\rightarrow$ $0.87\%$ IPC improvement over no $\mu$-op cache
$\rightarrow$ Increasing size of $\mu$-op cache does not help



*64Kops*
*32Kops*
*16Kops*
*8Kops*   *0.17%*

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

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

$\rightarrow$ Server workloads overwhelm current $\mu$-op caches
  $\rightarrow$ $0.87\%$ IPC improvement over no $\mu$-op cache
$\rightarrow$ Increasing size of $\mu$-op cache does not help



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

$\rightarrow$ Server workloads overwhelm current $\mu$-op caches
  $\rightarrow$ $0.87\%$ IPC improvement over no $\mu$-op cache
$\rightarrow$ Increasing size of $\mu$-op cache does not help



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

$\rightarrow$ Server workloads overwhelm current $\mu$-op caches
  $\rightarrow$ $0.87\%$ IPC improvement over no $\mu$-op cache
$\rightarrow$ Increasing size of $\mu$-op cache does not help



| | |
|---|---|
| 64Kops | 1.18% |
| 32Kops | 0.88% |
| 16Kops | 0.54% |
| 8Kops | 0.17% |

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

→ Server workloads overwhelm current $\mu$-op caches
  → $0.87\%$ IPC improvement over no $\mu$-op cache
→ Increasing size of $\mu$-op cache does not help
→ Ideal $\mu$-op cache can provide $10.82\%$ improvement



| | |
|---|---|
| Ideal | 10.82% |
| 64Kops | 1.18% |
| 32Kops | 0.88% |
| 16Kops | 0.54% |
| 8Kops | 0.17% |

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

→ Server workloads overwhelm current $\mu$-op caches
  → $0.87\%$ IPC improvement over no $\mu$-op cache
→ Increasing size of $\mu$-op cache does not help
→ Ideal $\mu$-op cache can provide $10.82\%$ improvement



| | |
|---|---|
| Ideal | 10.82% |
| 64Kops | 1.18% |
| 32Kops | 0.88% |
| 16Kops | 0.54% |
| 8Kops | 0.17% |

IPC Impro

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

$\rightarrow$ The FTQ is unable to hide the L1I miss latency on branch misprediction

$\rightarrow$ The FTQ is unable to hide the L1I miss latency on branch misprediction

- The FTQ is cleared on a branch misprediction

$\rightarrow$ The FTQ is unable to hide the L1I miss latency on branch misprediction

- The FTQ is cleared on a branch misprediction

*Alternate Path μ-op Cache Prefetching @ISCA'51*
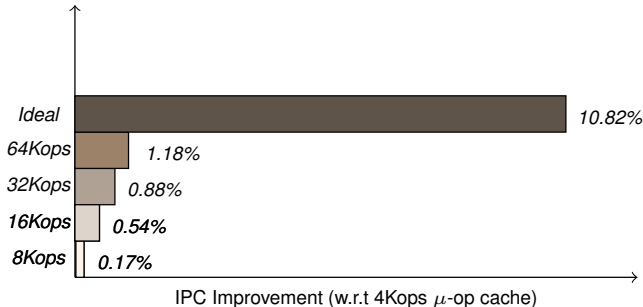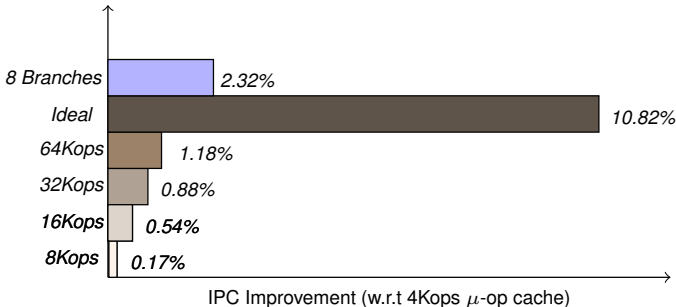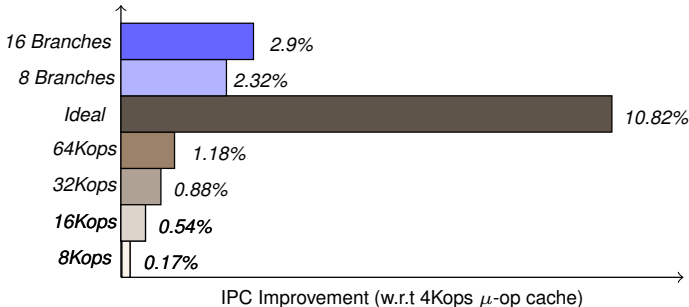
$\rightarrow$ The FTQ is unable to hide the L1I miss latency on branch misprediction

- The FTQ is cleared on a branch misprediction

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*
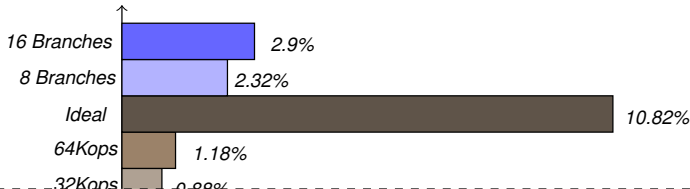
$\rightarrow$ The FTQ is unable to hide the L1I miss latency on branch misprediction

- The FTQ is cleared on a branch misprediction

$\rightarrow$ What if the correct path was always in the $\mu$-op cache after a pipeline flush due to branch misprediction?

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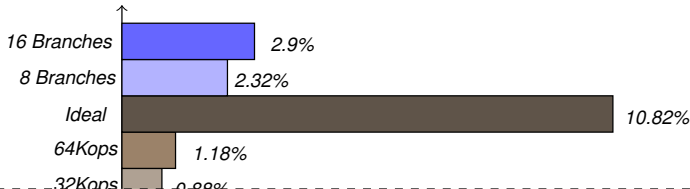→ What if the correct path was always in the $\mu$-op cache after a pipeline flush due to branch misprediction?



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

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

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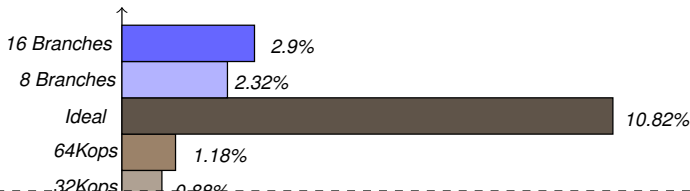→ What if the correct path was always in the $\mu$-op cache after a pipeline flush due to branch misprediction?



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

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

$\rightarrow$ The FTQ is unable to hide the L1I miss latency on branch misprediction

- The FTQ is cleared on a branch misprediction

$\rightarrow$ What if the correct path was always in the $\mu$-op cache after a pipeline flush due to branch misprediction?



| | |
|---|---|
| 16 Branches | 2.9% |
| 8 Branches | 2.32% |
| Ideal | 10.82% |
| 64Kops | 1.18% |
| 32Kops | 0.88% |
| 16Kops | 0.54% |
| 8Kops | 0.17% |

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

$\rightarrow$ The FTQ is unable to hide the L1I miss latency on branch misprediction

- The FTQ is cleared on a branch misprediction

$\rightarrow$ What if the correct path was always in the $\mu$-op cache after a pipeline flush due to branch misprediction?



| | |
|---|---|
| 16 Branches | 2.9% |
| 8 Branches | 2.32% |
| Ideal | 10.82% |
| 64Kops | 1.18% |
| 32Kops | 0.88% |

**Key Insight:**

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

$\rightarrow$ The FTQ is unable to hide the L1I miss latency on branch misprediction

- The FTQ is cleared on a branch misprediction

$\rightarrow$ What if the correct path was always in the $\mu$-op cache after a pipeline flush due to branch misprediction?



| | |
|---|---|
| 16 Branches | 2.9% |
| 8 Branches | 2.32% |
| Ideal | 10.82% |
| 64Kops | 1.18% |
| 32Kops | 0.88% |

**Key Insight:**

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

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

$\rightarrow$ The FTQ is unable to hide the L1I miss latency on branch misprediction
  - The FTQ is cleared on a branch misprediction

$\rightarrow$ What if the correct path was always in the $\mu$-op cache after a pipeline flush due to branch misprediction?



| | |
|---|---|
| 16 Branches | 2.9% |
| 8 Branches | 2.32% |
| Ideal | 10.82% |
| 64Kops | 1.18% |
| 32Kops | 0.88% |

**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

Ⓘ Identifies a hard-to-predict conditional branch (H2P)



H2P Branch ────────────────→ Predicted Path

*Alternate Path μ-op Cache Prefetching @ISCA'51*

① Identifies a hard-to-predict conditional branch (H2P)
② Generate addresses on alternate path

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

② Generate addresses on alternate path

③ Prefetch the alternate path to the $\mu$-op cache

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

① Identifies a hard-to-predict conditional branch (H2P)
② Generate addresses on alternate path
③ Prefetch the alternate path to the $\mu$-op cache



Alternate Path

*Prefetch to $\mu$-op cache*

Predicted Path

H2P Branch

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

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

$\rightarrow$ **H2P Branch:** a branch which has high chance of being mispredicted

→ H2P Branch: a branch which has high chance of being mispredicted
  → TAGE-Conf[2]

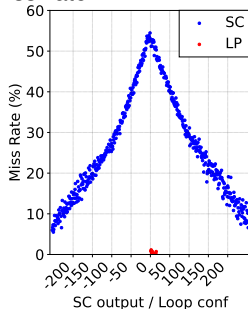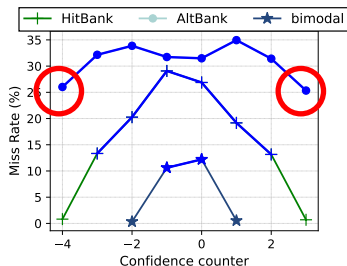*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

→ **H2P Branch:** a branch which has high chance of being mispredicted
  → TAGE-Conf[2]
    ● **Not saturated predictions** from AltBank, HitBank & BiModal

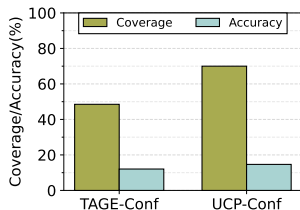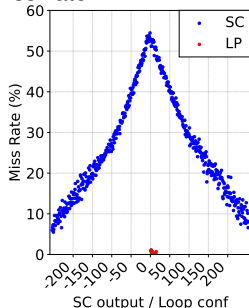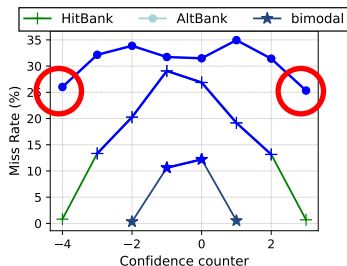*Alternate Path μ-op Cache Prefetching @ISCA'51*

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

    → TAGE-Conf[2]

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

    → UCP-Conf



[2]Seznec et. al. *Storage free confidence estimation for the TAGE branch predictor*

*Alternate Path μ-op Cache Prefetching @ISCA'51*

→ **H2P Branch:** a branch which has high chance of being mispredicted
  → TAGE-Conf[2]
    - Not saturated predictions from AltBank, HitBank & BiModal
  → UCP-Conf
    - All predictions from AltBanks shows high miss rate



[2]Seznec et. al. *Storage free confidence estimation for the TAGE branch predictor*

→ **H2P Branch:** a branch which has high chance of being mispredicted
  - → TAGE-Conf[2]
    - Not saturated predictions from AltBank, HitBank & BiModal
  - → UCP-Conf
    - All predictions from AltBanks shows high miss rate
    - SC shows high miss rate



[2]Seznec et. al. *Storage free confidence estimation for the TAGE branch predictor*

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

→ **H2P Branch:** a branch which has high chance of being mispredicted
  - → TAGE-Conf[2]
    - Not saturated predictions from AltBank, HitBank & BiModal
  - → UCP-Conf
    - All predictions from AltBanks shows high miss rate
    - SC shows high miss rate



[2]Seznec et. al. *Storage free confidence estimation for the TAGE branch predictor*

*Alternate Path μ-op Cache Prefetching @ISCA'51*

$\rightarrow$ Needs BPU

*Alternate Path μ-op Cache Prefetching @ISCA'51*

$\rightarrow$ Needs BPU
  - Banking

  - Add new predictors

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

$\rightarrow$ Needs BPU
- Banking
  - For server workloads BTB size is critical
  - Increase the number of banks from 16 to 32
- Add new predictors

*Alternate Path μ-op Cache Prefetching @ISCA'51*

$\rightarrow$ Needs BPU
- Banking
  - For server workloads BTB size is critical
  - Increase the number of banks from 16 to 32
- Add new predictors
  - Alt BP: 8KB TAGE-SC-L
  - Alt Indirect: 4KB ITTAGE
  - Alt RAS: 16-entry

$\rightarrow$ When to stop?

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

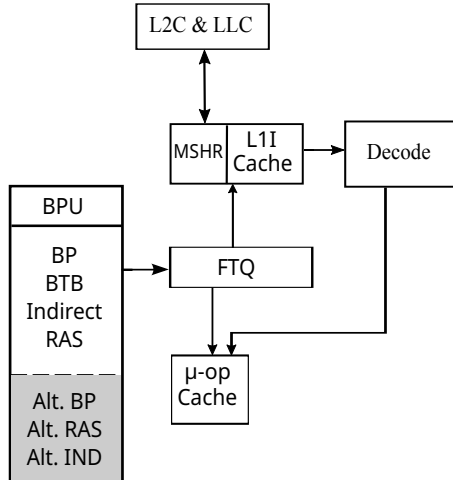$\rightarrow$ When to stop?

- Weight of each branch on the alternate path is accumulated
  - High confident $\rightarrow$ lower weight
  - Lower confident $\rightarrow$ higher weight



$AW = w_1 + w_2 + w_3$

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

$\rightarrow$ When to stop?

- Weight of each branch on the alternate path is accumulated
  - High confident $\rightarrow$ lower weight
  - Lower confident $\rightarrow$ higher weight
- BTB miss on the alternate path



AW = $w_1$ + $w_2$ + $w_3$

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

→ When to stop?

- Weight of each branch on the alternate path is accumulated
  - High confident → lower weight
  - Lower confident → higher weight
- BTB miss on the alternate path
- New H2P branch is detected



$AW = w_1 + w_2 + w_3$

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

# UCP
③ PREFETCHING

*Alternate Path μ-op Cache Prefetching @ISCA'51*

# UCP
## ③ PREFETCHING

*Alternate Path μ-op Cache Prefetching @ISCA'51*

UCP
③ PREFETCHING

L2C & LLC

MSHR | L1I Cache | Decode

BPU

Alt. FTQ ①

BP
BTB
Indirect
RAS

Alt. BP
Alt. RAS
Alt. IND

FTQ

μ-op Cache | MSHR

② ③

Tag Check

# UCP
③ PREFETCHING

*Alternate Path µ-op Cache Prefetching @ISCA'51*

# UCP
## ③ PREFETCHING

*Alternate Path μ-op Cache Prefetching @ISCA'51*

*Alternate Path μ-op Cache Prefetching @ISCA'51*

*Alternate Path μ-op Cache Prefetching @ISCA'51*

$\rightarrow$ ChampSim

$\rightarrow$ ChampSim

$\rightarrow$ Intel Alder Lake like microarchitecture

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

→ ChampSim

→ Intel Alder Lake like microarchitecture

→ Subset of CVP1[3] (traces showing $\geq$ 5% improvement with ideal $\mu$-op cache) [2 FP, 97 INT, 73 Crypto and 134 datacenter trace]

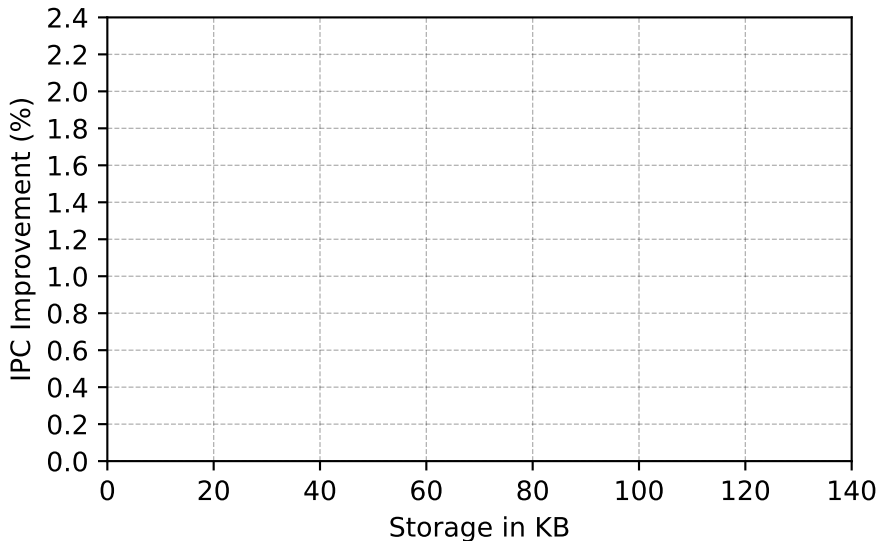[3]Feliu et. al. *Rebasing Microarchitectural Research with Industry Traces*

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

$\rightarrow$ ChampSim

$\rightarrow$ Intel Alder Lake like microarchitecture

$\rightarrow$ Subset of CVP1[3] (traces showing $\geq$ 5% improvement with ideal $\mu$-op cache) [2 FP, 97 INT, 73 Crypto and 134 datacenter trace]

$\rightarrow$ We execute 100M instructions, 50M warmup and 50M to collect stats

---

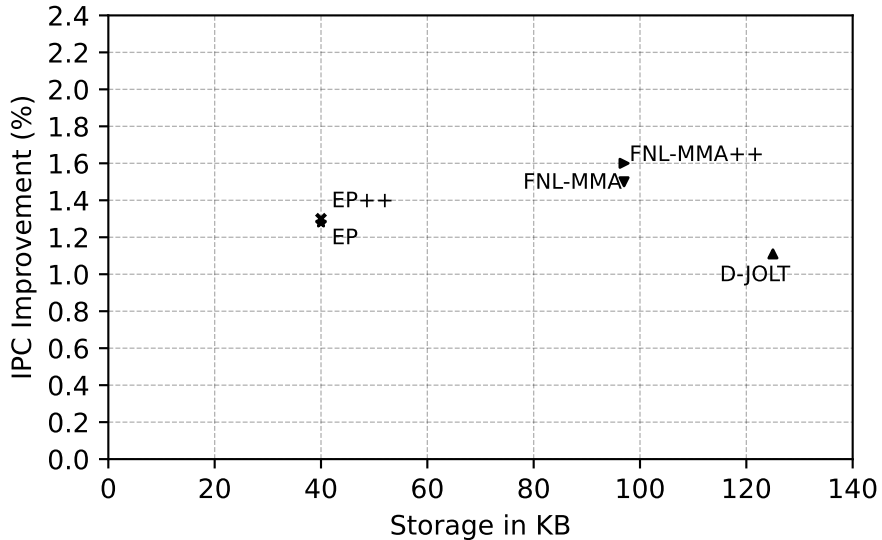[3]Feliu et. al. *Rebasing Microarchitectural Research with Industry Traces*

UCP sharing decoders with demand path, 1.8% improvement

Chart — IPC Improvement (%) vs Storage in KB:
- UCP-SharedDecoders (~13 KB, ~1.75%)
- UCP-L1I (~13 KB, ~1.65%)
- EP++ (~40 KB, ~1.3%)
- EP (~40 KB, ~1.3%)
- FNL-MMA (~97 KB, ~1.55%)
- FNL-MMA++ (~97 KB, ~1.65%)
- D-JOLT (~125 KB, ~1.1%)

Alternate Path $\mu$-op Cache Prefetching @ISCA'51

UCP with dedicated decoders, 2% improvement

Plot: IPC Improvement (%) vs Storage in KB

- UCP-AltDecoder
- UCP-SharedDecoders
- UCP-L1I
- EP++
- EP
- FNL-MMA
- FNL-MMA++
- D-JOLT

*Alternate Path μ-op Cache Prefetching @ISCA'51*

Ideal BTB scenario can give upto 2.2%

- • UCP-NoBTBConflict
- ▲ UCP-AltDecoder
- UCP-SharedDecoders
- UCP-L1I
- EP++
- EP
- FNL-MMA++
- ▶ FNL-MMA
- ▲ D-JOLT

IPC Improvement (%) vs Storage in KB

Scatter plot: IPC Improvement (%) vs Storage in KB

- UCP-NoBTBConflict
- UCP-AltDecoder
- UCP-SharedDecoders
- UCP-L1I
- TAGE-SC-Lx2
- FNL-MMA++
- FNL-MMA
- EP++
- EP
- D-JOLT

Doubling the TAGE-SC-L predictor (64KB → 128KB)

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

$\rightarrow$ $\mu$-op cache can be used for performance improvements

$\rightarrow$ $\mu$-op cache can be used for performance improvements

$\rightarrow$ We propose UCP ($\mu$-*op Cache Prefetching*)

$\rightarrow$ $\mu$-op cache can be used for performance improvements

$\rightarrow$ We propose UCP ($\mu$-op Cache Prefetching)

- Identifies a hard-to-predict branch

→ $\mu$-op cache can be used for performance improvements

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

  • Identifies a hard-to-predict branch

  • Prefetch critical instructions in the $\mu$-op cache

# ALTERNATE PATH $\mu$-OP CACHE PREFETCHING

**Sawan Singh**[1]    Arthur Perais[2]    Alexandra Jimborean[1]
Alberto Ros[1]

singh.sawan@um.es

Thank you for your attention!
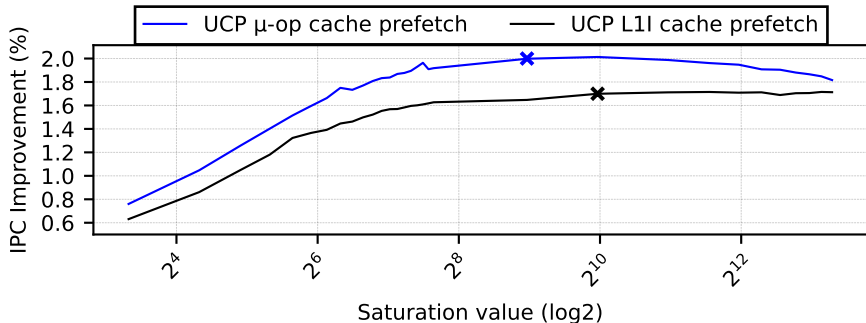
→ UCP reuses the BTB by doubling the number of BTB banks (from 16 to 32)

→ 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

*Alternate Path $\mu$-op Cache Prefetching @ISCA'51*

*Alternate Path μ-op Cache Prefetching @ISCA'51*