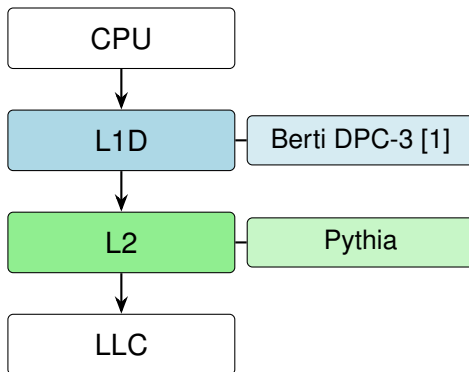# PUSHING THE LIMITS OF THE BERTI PREFETCHER

**Simranjit Singh**[1], Agustín Navarro-Torres[2], Alberto Ros[1]

[1]University of Murcia (simranjit.singh@um.es, aros@ditec.um.es)
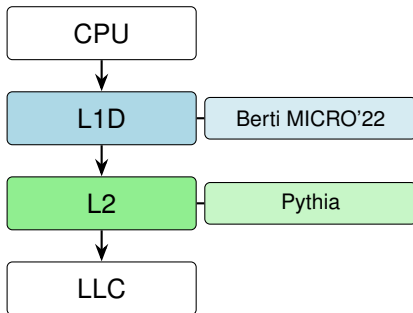[2]Universidad de Zaragoza (agusnt@unizar.es)

Feb 1, 2026

# BASELINE



CPU → L1D → L2 → LLC

L1D — Berti DPC-3 [1]

L2 — Pythia

[1] DPC-3 = 3rd Data Prefetching Championship (2019)

# BASELINE: BERTI DPC-3 → MICRO'22

```
┌──────────┐
│   CPU    │
└──────────┘
     │
     ▼
┌──────────┐      ┌──────────────────┐
│   L1D    │──────│  Berti MICRO'22  │
└──────────┘      └──────────────────┘
     │
     ▼
┌──────────┐      ┌──────────────────┐
│   L2     │──────│     Pythia       │
└──────────┘      └──────────────────┘
     │
     ▼
┌──────────┐
│   LLC    │
└──────────┘
```
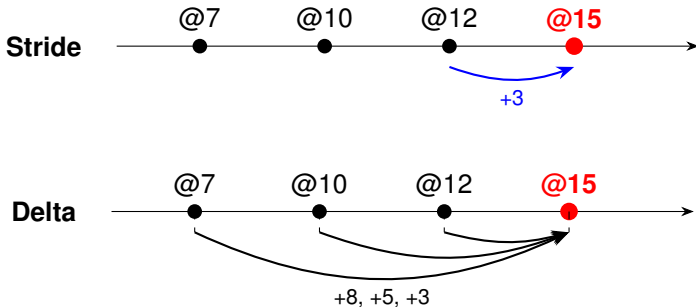
**MICRO'22 Improvements**

- Simplified design
- Virtual addresses
  (enables cross-page prefetching)
- Local per-IP deltas
  (vs memory regions)
- Timeliness via latency tracking
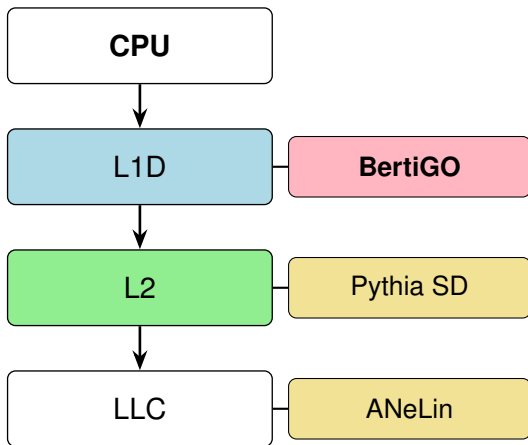- Confidence-based delta selection

**Stride vs Local Deltas**

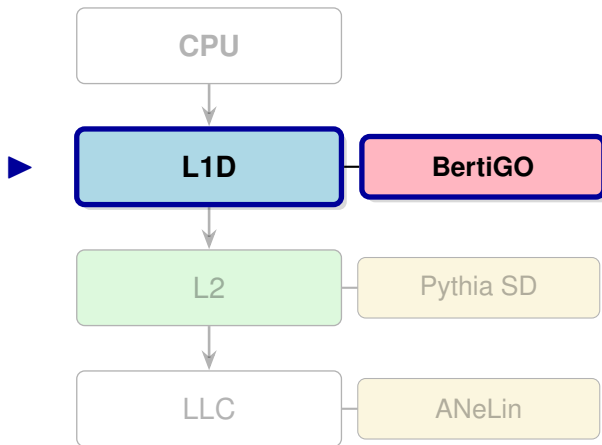Per-IP deltas • Latency tracking • Confidence-based L1D/L2 •
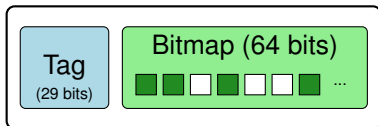Virtual addresses

# OUR APPROACH

Berti prefetches lines
**already in cache**

$\rightarrow$ Uses PQ slots
$\rightarrow$ Uses port slots

# SOLUTION 1: REGION-BASED BITMAP FILTER

**Filter Entry**

| Tag (29 bits) | Bitmap (64 bits) ■■□■□□□■ ... |
|---|---|

**How it works:**

1. On access/prefetch: **Set bit**
2. Before prefetch: **Check bit**
3. Bit set? → **Skip prefetch!**
4. On L1D eviction: **Clear bit**

**Learns Useless Prefetches**

If L2 prefetch never promoted to L1D:

→ No eviction
→ Bit stays set
→ Future requests blocked

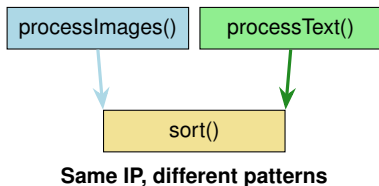**Storage:** 15.6 KB
1360 entries · NRU
Tracks 87,040 lines

Berti indexes by IP alone

**Can we add new input sources?**

# SOLUTION 2: CONTEXT-AWARE PREFETCHING

| processImages() | processText() |
|:---:|:---:|

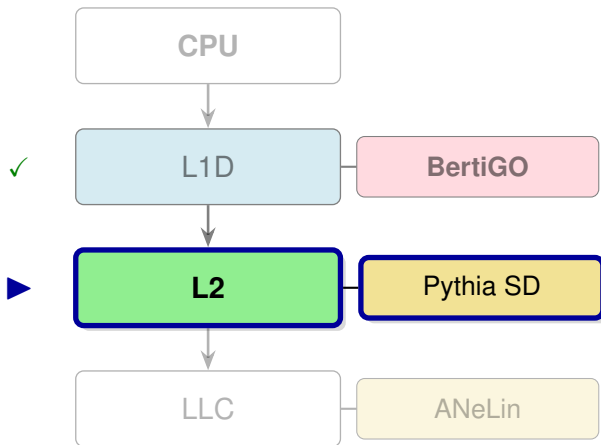| sort() |
|:---:|

**Same IP, different patterns**

**How it works:**

1. Track last 4 IPs
2. Hash into path signature
3. Query deltas using both IP alone and path signature
4. Merge predictions

**Highly effective for AI/ML**

Llama (full BW):
**+20-50%** vs IP-only

**Storage:** $< 0.01$ KB

# BERTIGO: L1D CONTRIBUTION



Full BW, over baseline

# OUR APPROACH

# PYTHIA: KEY IDEAS



**Prefetcher**

Features of memory request to address A (e.g., IP)

**Reward**

Prefetch from address A+offset (O)

**Processor & Memory Subsystem**

RL Learning • Multiple program context features • Bandwidth-aware rewards

# L2: SET-DUELING FOR PYTHIA

**2048 Cache Sets**



- 🟥 NoPref  🟦 IP  🟩 IP_Delta
- 🟧 IP∨IP_Delta  ⬛ IP∧IP_Delta

**How it works:**

1. Each set assigned to one candidate
2. Track miss rate (misses ÷ accesses)
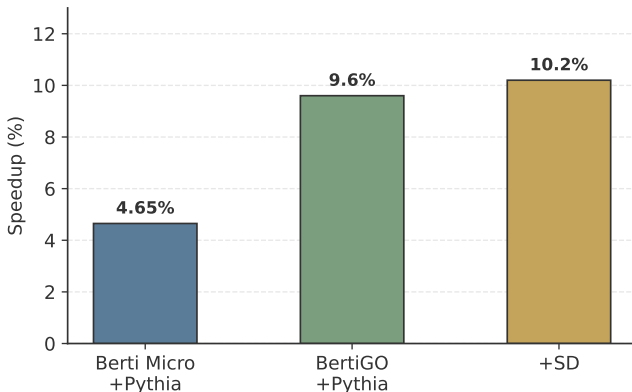3. Winner if ≥4% better miss rate than NoPrefetch

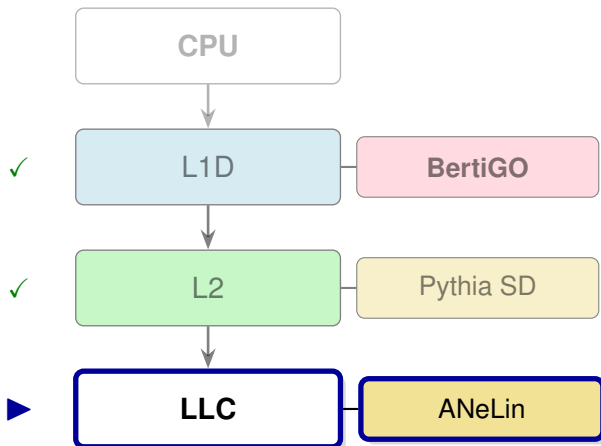**Workload-Driven**

Default policy not always optimal

→ Let workload decide

**Storage:** 90 bytes
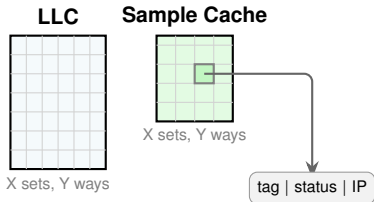
# RESULTS: BERTIGO + PYTHIA SD



Full BW, over baseline

# OUR APPROACH

# LLC: ANeLin (Adaptive Next-Line)

**LLC**    **Sample Cache**



X sets, Y ways

X sets, Y ways

tag | status | IP

**Two-Level Filtering**
Per-Core: enable for
workload?
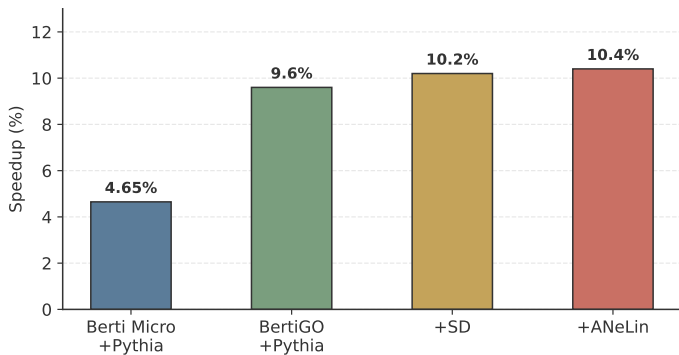Per-IP: enable for this IP?

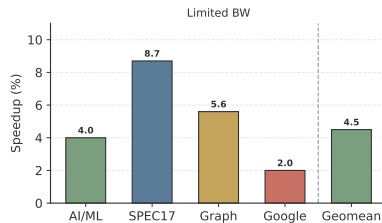**Storage:** 209 KB

[1] Timely: $>$ mean lat $\rightarrow$ 1 useful
[2] Late: $> 0.2\times$ mean lat $\rightarrow$ 0.5 useful

**How it works:**

1. On access: track next-line in sample cache; prefetch if enabled
2. Demand hit $\rightarrow$ timely[1] or late[2]
3. Unused eviction $\rightarrow$ useless
4. Saturate: enable next-line if useful $\gg$ useless; decay counters 75%

# RESULTS: FULL CONFIGURATION



Full BW, over baseline

# EVALUATION

# CONCLUSION

**L1D: BertiGO**

> **Region Filter**
> • Eliminates redundant requests
> • Learns useless L2 prefetches

> **IP-Path Signatures**
> • Context-aware predictions
> • Effective for AI/ML

**L2: Pythia SD**

> **Set-Dueling**
> • First applied to prefetcher feature selection
> • Choose optimal policy per workload

**LLC: ANeLin**

> **Adaptive Next-Line**
> • Per-IP and per-core learning

# Thank You

Questions?
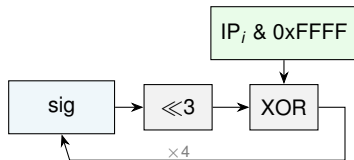
L1D: BertiGO | L2: Pythia Set-Dueling | LLC: ANeLin

Simranjit Singh, Agustín Navarro-Torres, Alberto Ros — DPC-4

# BACKUP: IP-PATH SIGNATURE

```
sig ← 0
for i = 3 downto 0:
  sig ← (sig ≪ 3) ⊕
        (IPᵢ & 0xFFFF)
return sig
```



**Context A**

0x12, 0x34, 0x56, 0x78, 0xAB

⇓

0x15EEB

**Context B**

0xFF, 0x11, 0x22, 0x33, 0xAB

⇓

0xFDBB3