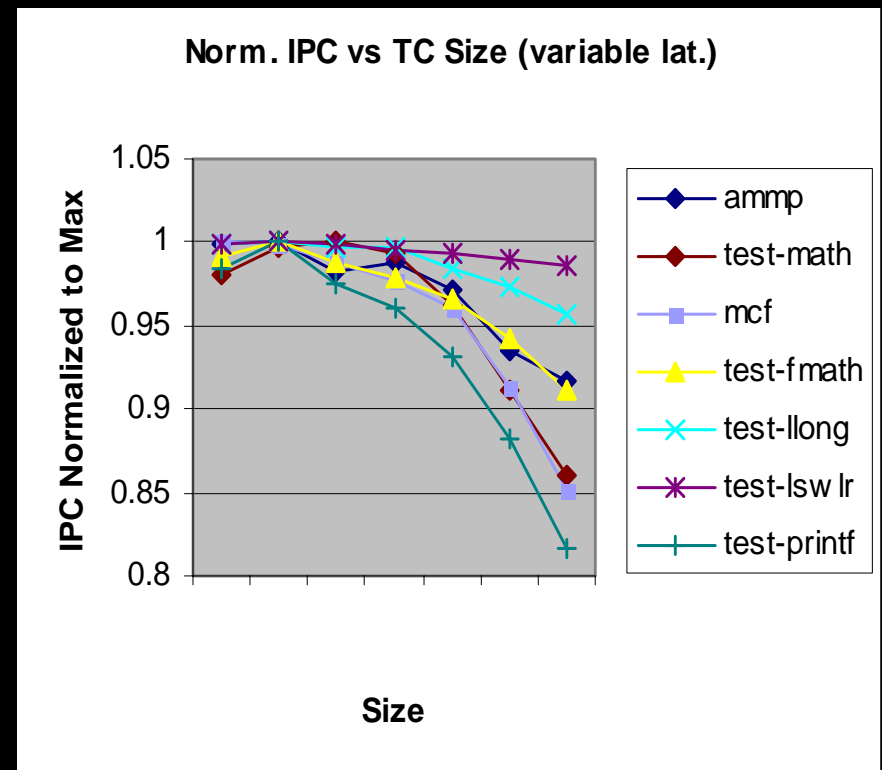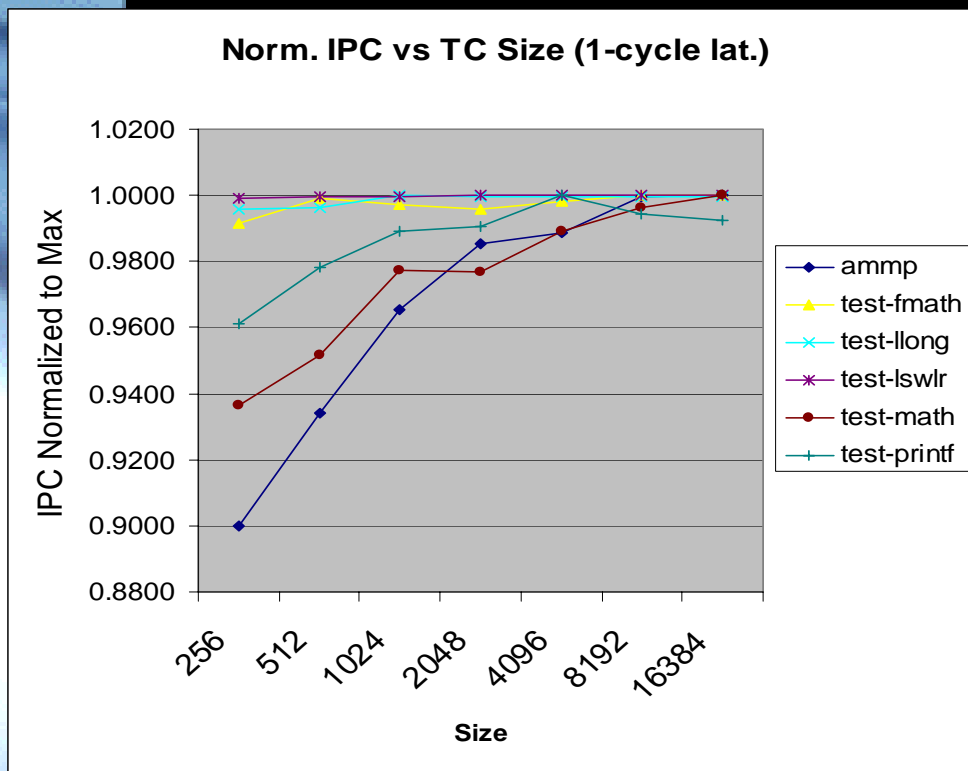# *Frequently Used Trace Cache*

```
[efurbish@hoss simplesim-2.0]> wc -l trace_cache.c
    666
```

Advanced Computer Architecture,
under Dr. Scott Rixner

Marc Power

Eric Furbish

Indraneel Datta

# M

- Large trace caches suffer increased latency

  - Lowers performance below ideal



**Norm. IPC vs TC Size (1-cycle lat.)**

IPC Normalized to Max

1.0200
1.0000
0.9800
0.9600
0.9400
0.9200
0.9000
0.8800

256  512  1024  2048  4096  8192  16384

Size

- ammp
- test-fmath
- test-llong
- test-lswlr
- test-math
- test-printf



**Norm. IPC vs TC Size (variable lat.)**

IPC Normalized to Max

1.05
1
0.95
0.9
0.85
0.8

Size

- ammp
- test-math
- mcf
- test-fmath
- test-llong
- test-lswlr
- test-printf

- We'd like to minimize the performance lost to latency.

# C

- To recover from the ill effects of latency in the trace cache, we need some *low latency* means of getting at the same data.

- Normally, solve by implementing L1/L2 cache hierarchy

  - Contention in the L1 can evict frequently used lines with infrequently used ones

  - Contention could possibly be reduced by intelligently filling the L1

- *Frequently Used Trace Cache (FUTC) –*

  - Single-cycle latency, small L1

  - Judiciously filled from L2 with frequently used lines

- Given the new gizmo, we hoped:

  - $IPC(TC_s + FUTC_i) > IPC(TC_s + L1TC_i) > IPC(TC_s)$

# FU

- One saturating counter per trace cache line

    - Incremented on read hit

- Counters cleared on TC writes

    - Rewrites not propagated into FUTC

        - Saves read port

        - Maintains logical meaning of counters (apply to a *specific* trace)

- Lines with counters over threshold are promoted

    - Promoted on a read hit

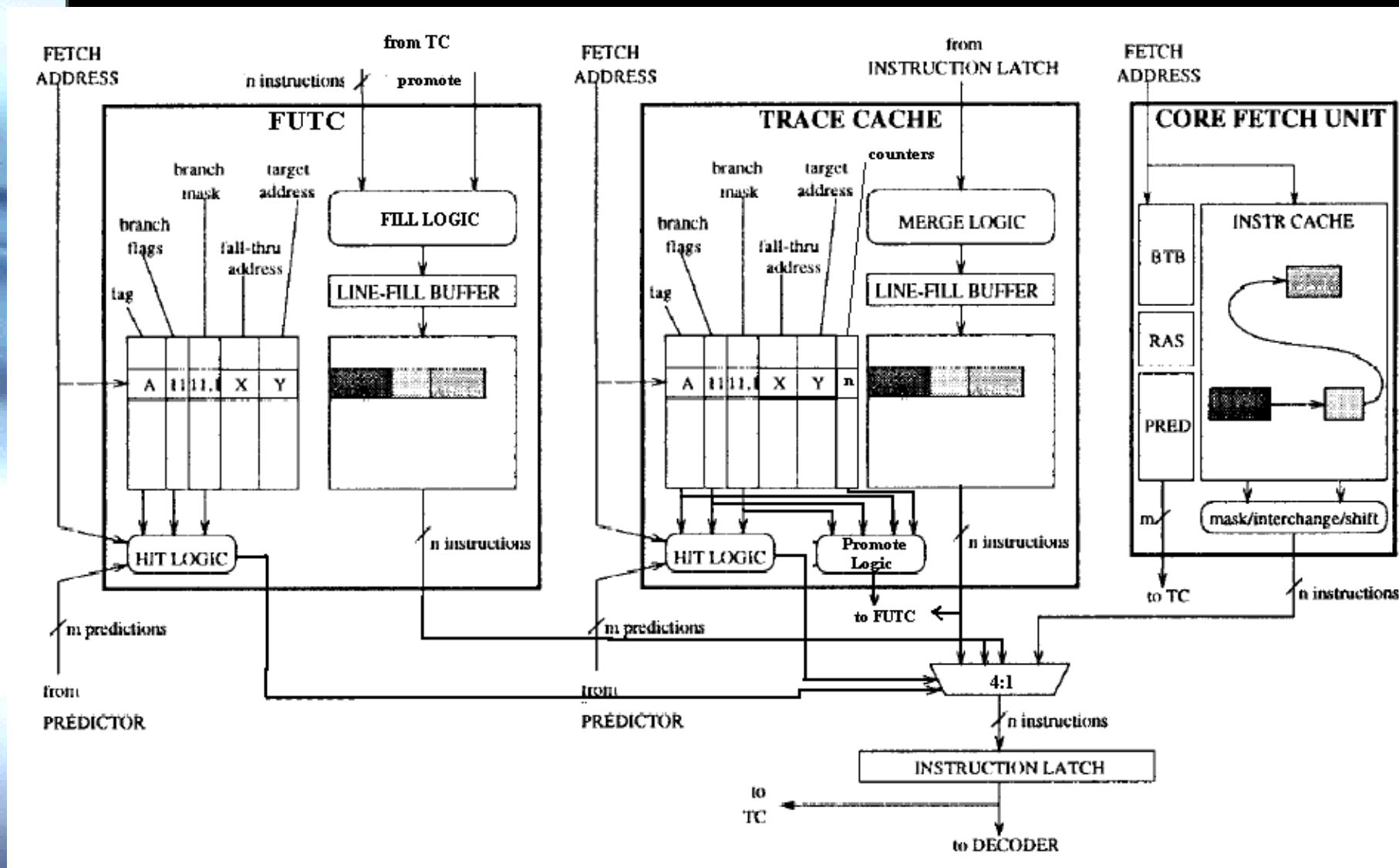    - Counter is not reset

# FUTC



Figure Modified from (E. Rotenbert, S. Bennett, J. Smith. Trace cache: a low latency approach to high bandwidth instruction fetching. Tech Report 1310, CS Dept., Univ. of Wisc.-Madison, 1996)

# A...

- Size of FUTC (8-64KB is reasonable) in order to remain under single-cycle latency

- Counters consume 1-3KB for 1MB trace cache

  - Assumes threshold of 1-8 (1-3 bits)

- Power cost is minimal

  - No worse than L1 trace cache

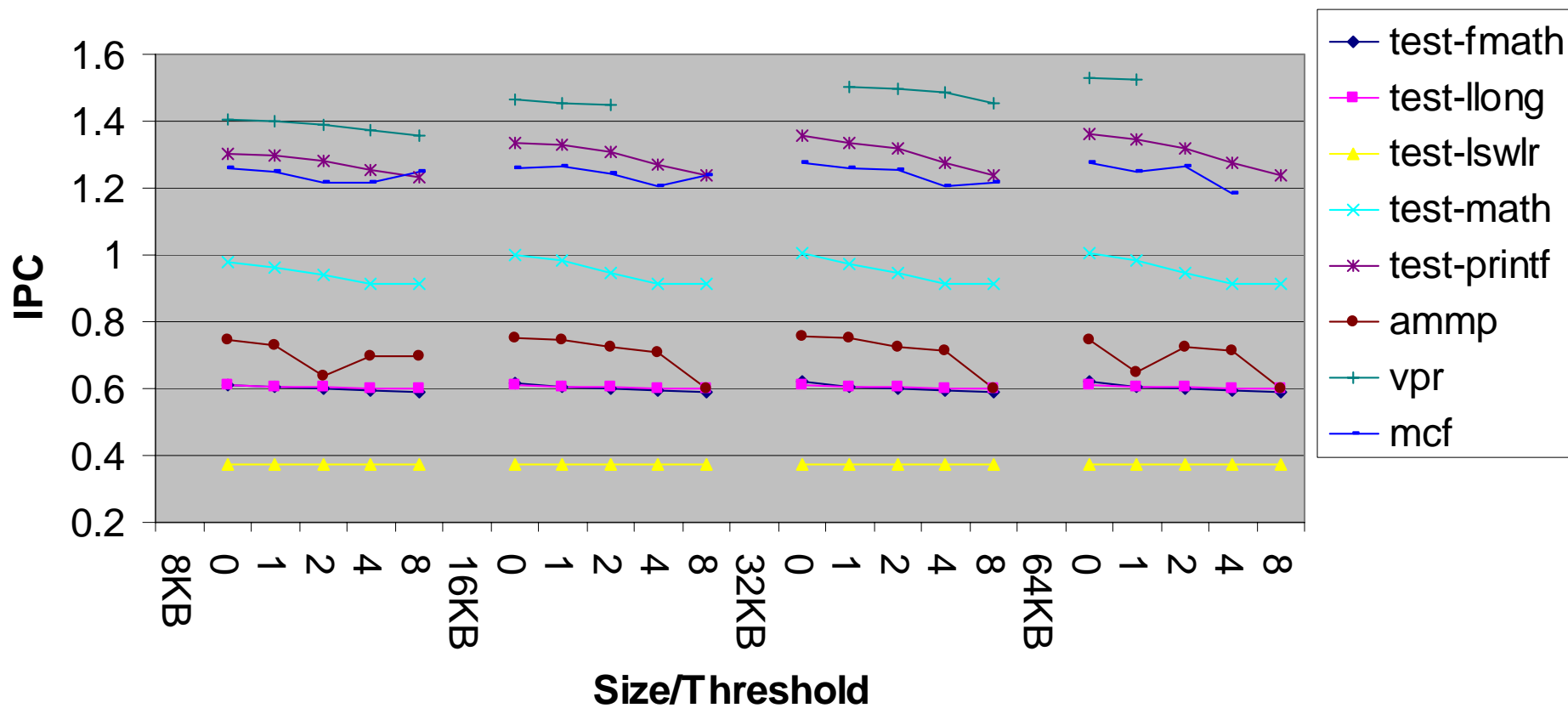  - Single-ported for speed and power

M

1. Use "infinite" backend to stress fetch mechanism

2. Select 'Optimal' TC Size

   - Based on ideal (single-cycle) performance vs. performance with latency accounted for
   - Pick best candidate for improvement with FUTC

3. Run varying "approximate L1" sizes (8-64KB)

   - L1 = FUTC with threshold 0

4. Run varying FUTC sizes (8-64KB) with thresholds of 1, 2, 4 and 8

5. Compare TC/FUTC and L1/FUTC

# Ve...

- Branch predictor accuracies of 90-99%
  - Exactly the range we would expect
- Performance and trace cache hit rates very similar to trace cache paper
- Dispatch and commit streams are identical to unmodified sim-outorder
  - Program flow is guaranteed to be correct
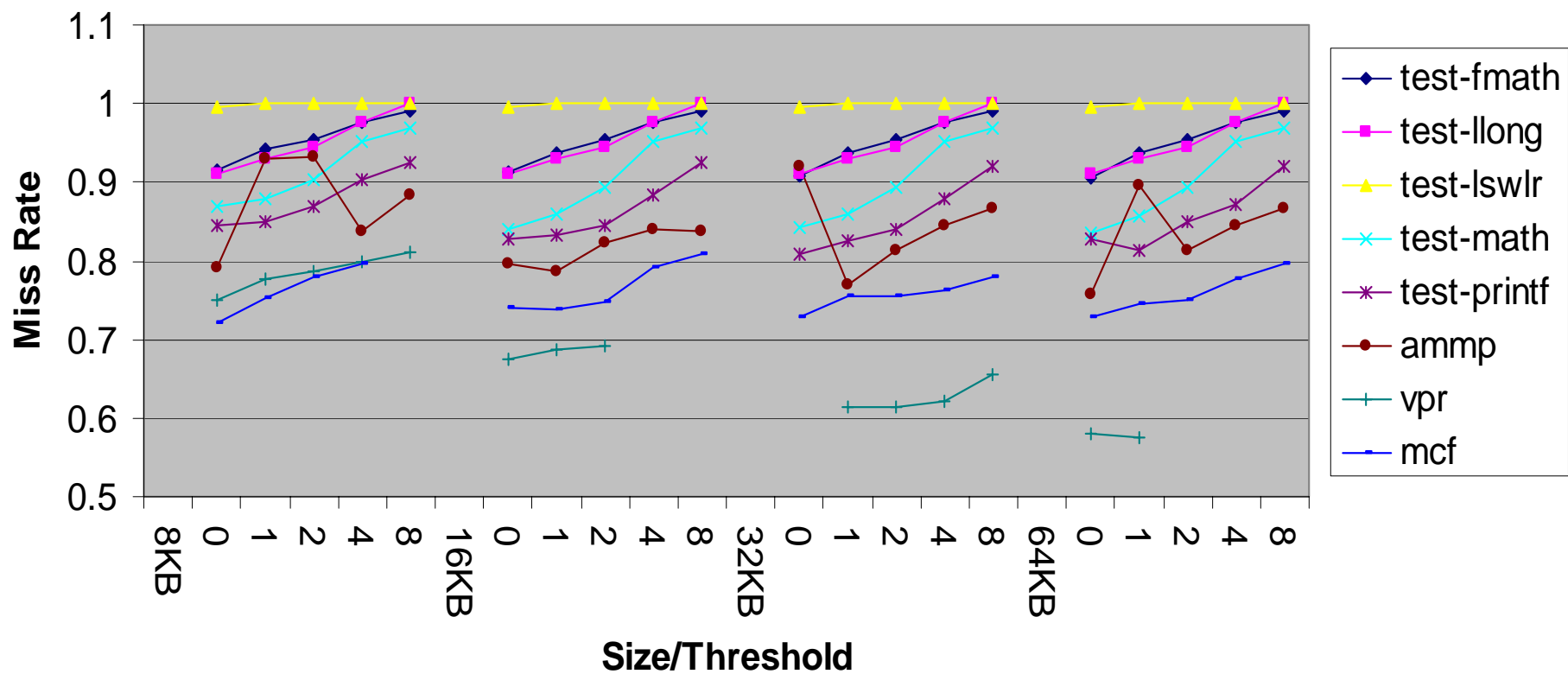- It has seen billions of instructions without flaws

# Re



**IPC vs. FUTC size/Threshold**
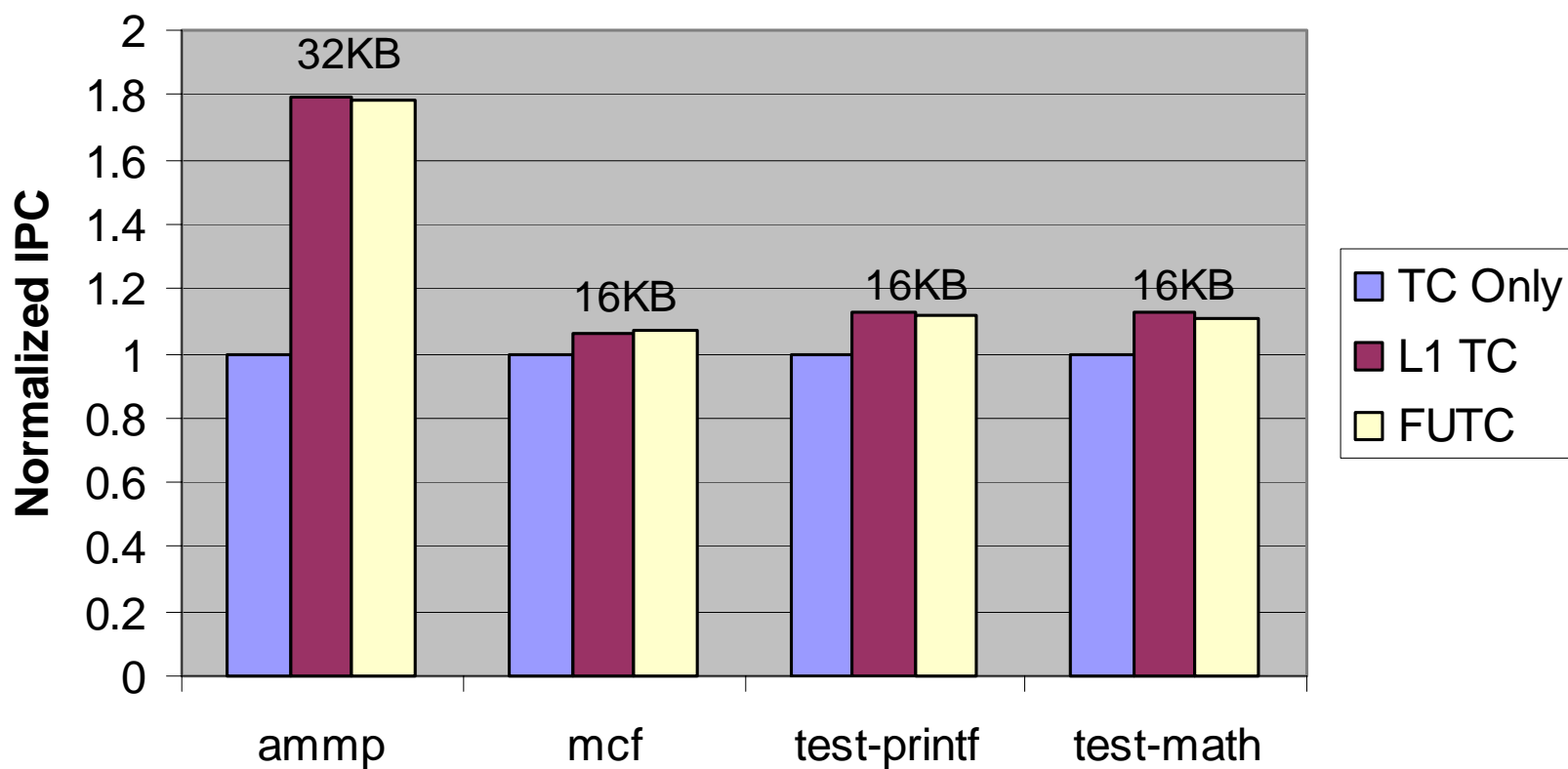
Legend:
- test-fmath
- test-llong
- test-lswlr
- test-math
- test-printf
- ammp
- vpr
- mcf

Y-axis: IPC (0.2, 0.4, 0.6, 0.8, 1, 1.2, 1.4, 1.6)

X-axis: Size/Threshold (8KB: 0 1 2 4 8, 16KB: 0 1 2 4 8, 32KB: 0 1 2 4 8, 64KB: 0 1 2 4 8)

# Re



**FUTC Miss Rate vs FUTC Size/Threshold**

Legend: test-fmath, test-llong, test-lswlr, test-math, test-printf, ammp, vpr, mcf

Y-axis: Miss Rate (0.5 to 1.1)

X-axis: Size/Threshold (8KB, 16KB, 32KB, 64KB with values 0, 1, 2, 4, 8)

# Be

**Best-Case FUTC vs L1 and TC**



Bar chart showing Normalized IPC for benchmarks ammp (32KB), mcf (16KB), test-printf (16KB), and test-math (16KB), comparing TC Only, L1 TC, and FUTC.

# Fu

- Run larger programs
  - Increased contention could make FUTC effective

- Use partial matching and inactive issue
  - Higher TC hit rate = more promotions = more contention

# C

- Hypothesis was partly correct
  - Small auxiliary cache always helps over TC
- Contention in L1 outweighed by FUTC "warmup"
- New Hypothesis:
  - $IPC(TC_s + L1\ TC_i) > IPC(TC_s + FUTC_i) > IPC(TC_s)$
  - Really: $IPC(TC_s + AUX\text{-}TC_i) > IPC(TC_s)$