# LTRF: Enabling High-Capacity Register Files for GPUs via Hardware/Software Cooperative Register Prefetching

Mohammad Sadrosadati

Seyed Borna Ehsani

Mario Drumond

Rachata Ausavarungnirun

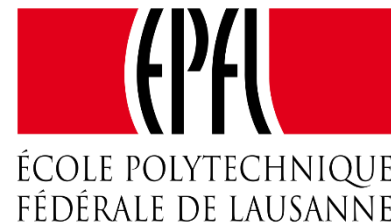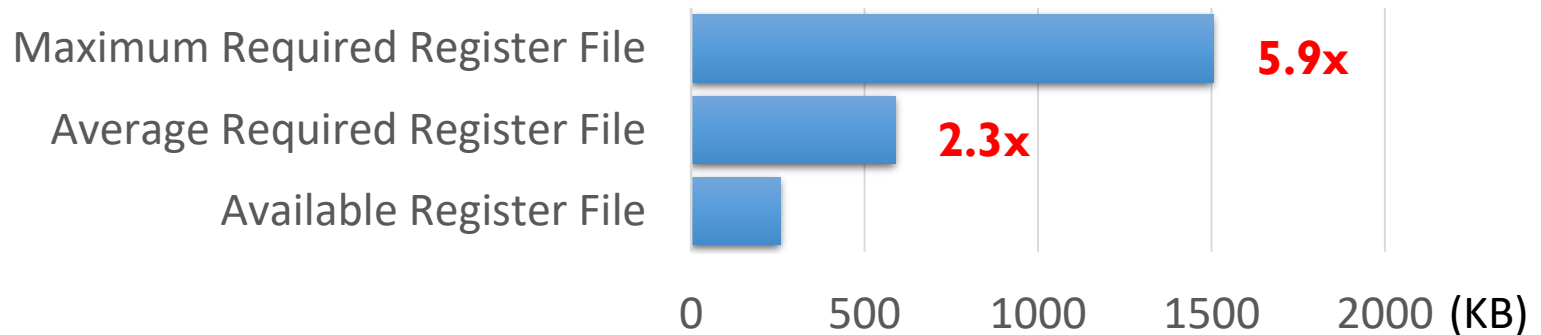**Amirhossein Mirhosseini**

Hamid Sarbazi-Azad

Babak Falsafi

Onur Mutlu

# Register file size limits GPU scalability

- Register file already accounts for 60% of on-chip storage

- But, there is still demand for more registers to achieve maximum performance and concurrency

Maximum Required Register File **5.9x**

Average Required Register File **2.3x**

Available Register File
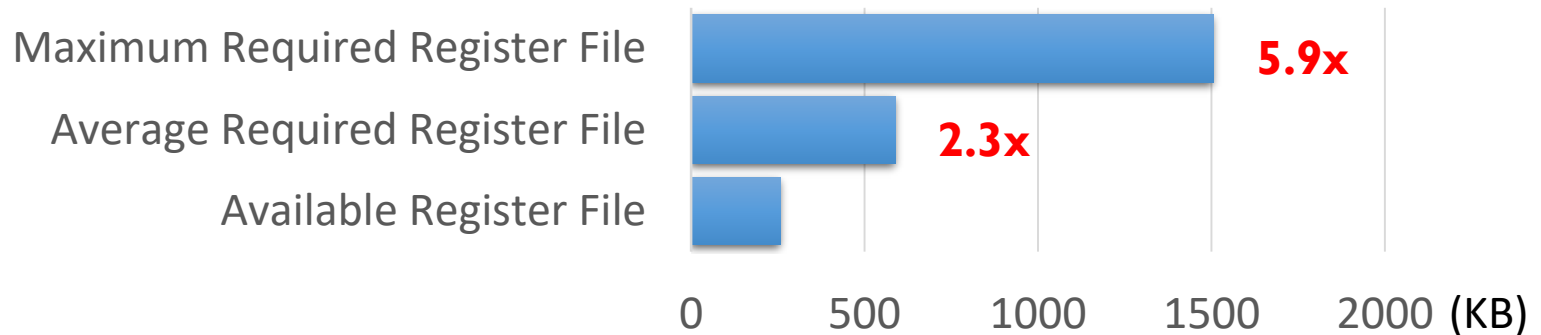
0    500    1000    1500    2000 (KB)
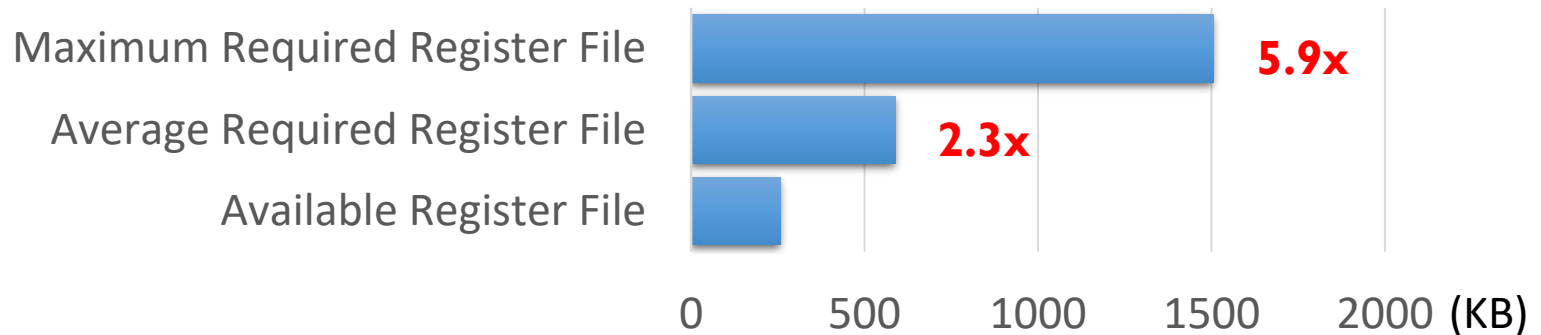
# Register file size limits GPU scalability

- Register file already accounts for 60% of on-chip storage

- But, there is still demand for more registers to achieve maximum performance and concurrency

Maximum Required Register File **5.9x**

Average Required Register File **2.3x**

Available Register File

0    500    1000    1500    2000 (KB)

•Unfortunately,  all mechanisms to expand RF capacity without large area/power overheads significantly increase access latency

- Emerging technologies, register file compression/virtualization, etc.

3

# Register file size limits GPU scalability

- Register file already accounts for 60% of on-chip storage

- But, there is still demand for more registers to achieve maximum performance and concurrency



Maximum Required Register File — 5.9x
Average Required Register File — 2.3x
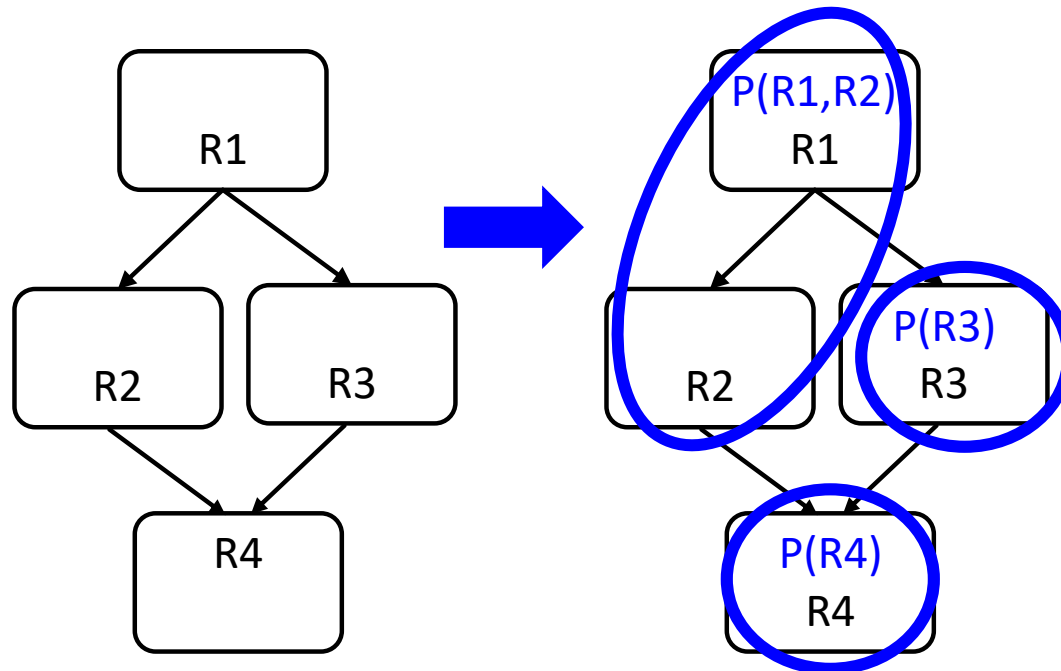Available Register File

0    500    1000    1500    2000  (KB)

- Unfortunately,  all mechanisms to expand RF capacity without large area/power overheads significantly increase access latency
  - Emerging technologies, register file compression/virtualization, etc.

## Goal: Tolerate register file latencies

# Contributions (1)

- **Compiler-driven Register Prefetching**
  - Break control flow graph into "prefetch subgraphs"
  - Prefetch registers at the beginning of each subgraph
  - Interval analysis to identify prefetch subgraphs

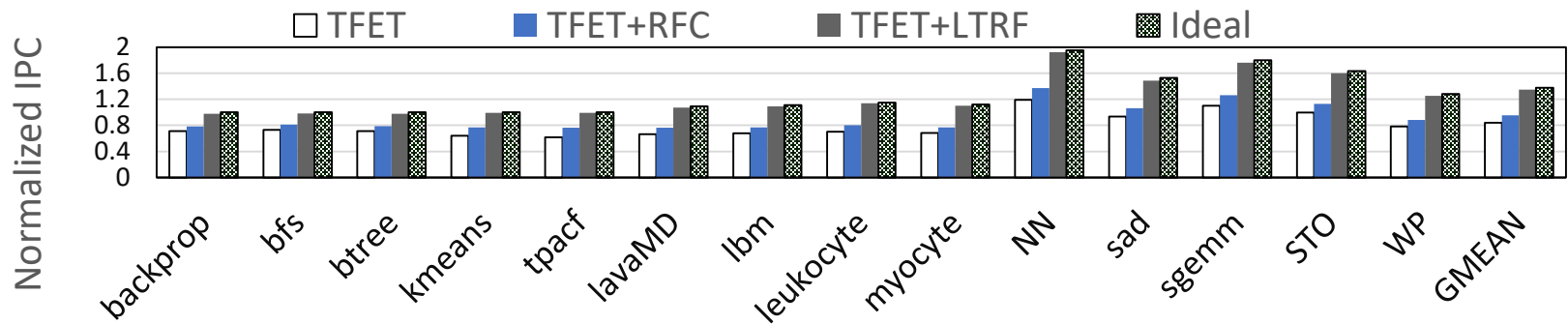# Contributions (2)

- **Latency Tolerant Register File (LTRF)**
    - "2-level" main register file + register cache
    - Performs prefetch ops while executing other warps
    - Tolerates up to 6x higher RF access latencies
    - Paves the way for various area/power optimizations

# Contributions (2)

- **Latency Tolerant Register File (LTRF)**
  - "2-level" main register file + register cache
  - Performs prefetch ops while executing other warps
  - Tolerates up to 6x higher RF access latencies
  - Paves the way for various area/power optimizations



Example LTRF deployment enables 8× larger register file and 32% higher performance

# LTRF: Enabling High-Capacity Register Files for GPUs via Hardware/Software Cooperative Register Prefetching
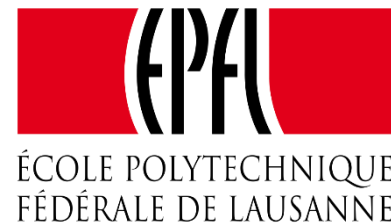
Mohammad Sadrosadati

Seyed Borna Ehsani

Mario Drumond

Rachata Ausavarungnirun

**Amirhossein Mirhosseini**

Hamid Sarbazi-Azad

Babak Falsafi

Onur Mutlu

# LTRF: Enabling High-Capacity Register Files for GPUs via Hardware/Software Cooperative Register Prefetching

Mohammad Sadrosadati

Seyed Borna Ehsani

Mario Drumond

Rachata Ausavarungnirun

**Amirhossein Mirhosseini**

Hamid Sarbazi-Azad

Babak Falsafi

Onur Mutlu

Please attend our talk at 2:00 in Virginia EF