# mlin: Rethinking and Rebooting `gprof` for the Multicore Era

**nino Garcia**, Donghwan Jeon, Chris Louie, Michael B. Taylor

outer Science & Engineering Department
rsity of California, San Diego

# tivating a "`gprof` for parallelization"

*w effective are programmers at picking the right parts of a gram to parallelize?*

- User study* we performed at UC San Diego (UCSD IRB #100056)

- First and second year CS graduate students

- Users parallelize their programs and submit to job queue for timing

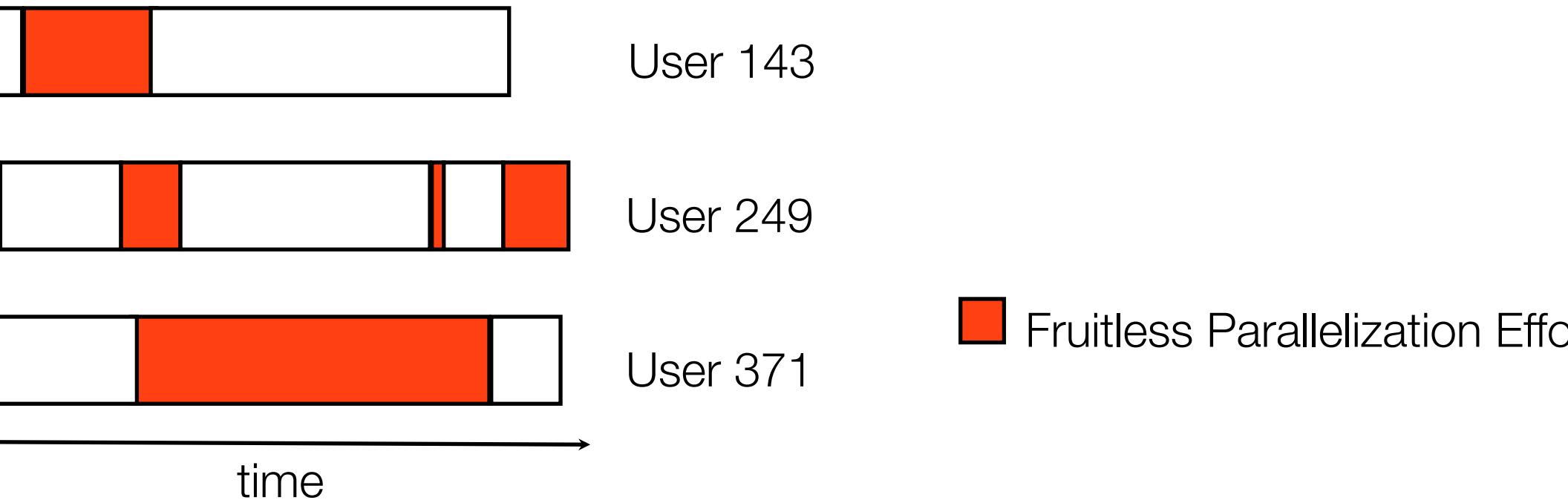- 32-core AMD machine, Cilk++, access to gprof

- Students were graded based on effectiveness of their parallel spee tudents told serial optimization would not help their grade

amined student's activities to determine result of efforts



User 143

User 249

User 371

□ Fruitless Parallelization Eff

time

nificant fraction of fruitless effort because of three basic probl

*Low Parallelism*: Region was not parallel enough

*Low Coverage*: Region's execution time was too small

**gprof** answers the question:

*"What parts of this program should I spend time **optimizing**?"*

**Kremlin** answers the question:

*"What parts of this program should I spend time **parallelizing**?"*

# mlin's Usage Model

age model inspired by gprof

ake CC=kremlin-cc

/tracking lolcats

remlin tracking --personality=openmp

1. Compile instrumented binary

2. Profile with sample input

3. Run analysis tool to create

| ile (lines) | Self-P | Cov (%) |
|---|---|---|
| mageBlur.c (49-58) | 145.3 | 9.7 |
| mageBlur.c (37-45) | 145.3 | 8.7 |
| etInterpPatch.c (26-35) | 25.3 | 8.9 |
| alcSobel_dX.c (59-68) | 126.2 | 8.1 |
| alcSobel_dX.c (46-55) | 126.2 | 8.1 |

*• Hierarchical Critical Path Analysis* (HCPA)

▸ Quantifies self parallelism in each program region

**Parallelism Discovery**
*"What's the potential parallel speedup of
each part of this program?"*

*• Self-Parallelism*

▸ Estimates ideal parallel speedup of a specific region

*• Planning Personalities*

**Parallelism Planning**
*"What regions must I parallelize to get
the maximum benefit on this system?"*

▸ Incorporates target specific constraints in parallelization

Parallelization

sting Technique: 1980's-era Critical Path Analysis (CPA)

Finds critical path through the dynamic execution of a program

Mainly used in research studies to quantify limits of parallelism

$$\text{parallelism} = \frac{\text{work}}{\text{critical path length}}$$

**work ~= # of instrs**

○ instruction

↘ data or control dependence

**critical path (cp)**

uates program's potential for parallelization under relatively
nistic assumptions

oser approximation to what human experts can achieve

rsus pessimistic static analysis in automatic parallelizing compilers

It is predictive of parallelism after typical parallelization
formations

., Loop interchange, loop fission, locality enhancement

A is typically run on an entire program

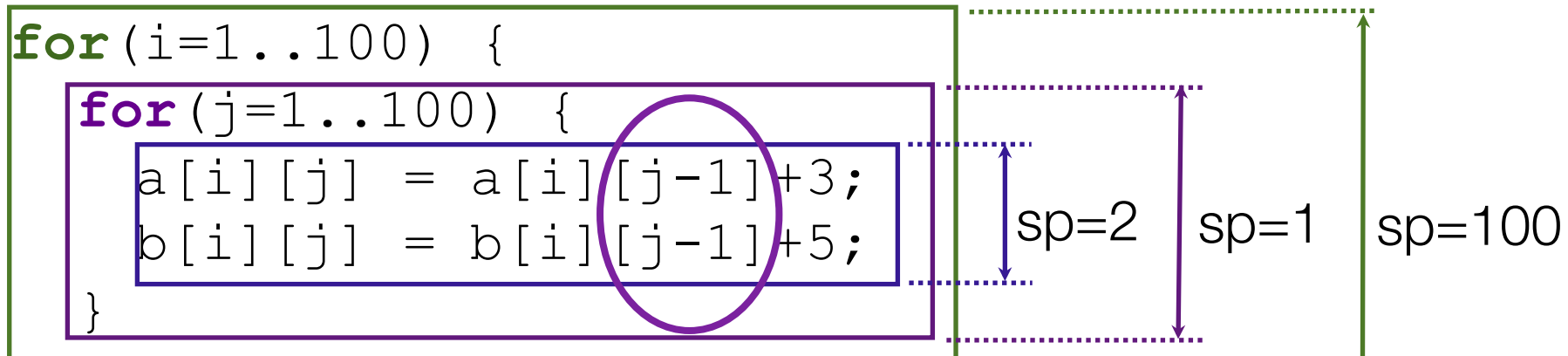Not helpful for identifying specific regions to parallelize

Doesn't help evaluate execution time of a program if only a subset he program is parallelized

rarchical CPA is a region-based analysis

Self-Parallelism (sp) identifies parallelism in specific regions

Provides basis for estimating parallel speedup of individual regions

```
for(i=1..100) {
  for(j=1..100) {
    a[i][j] = a[i][j-1]+3;
    b[i][j] = b[i][j-1]+5;
  }
}
```
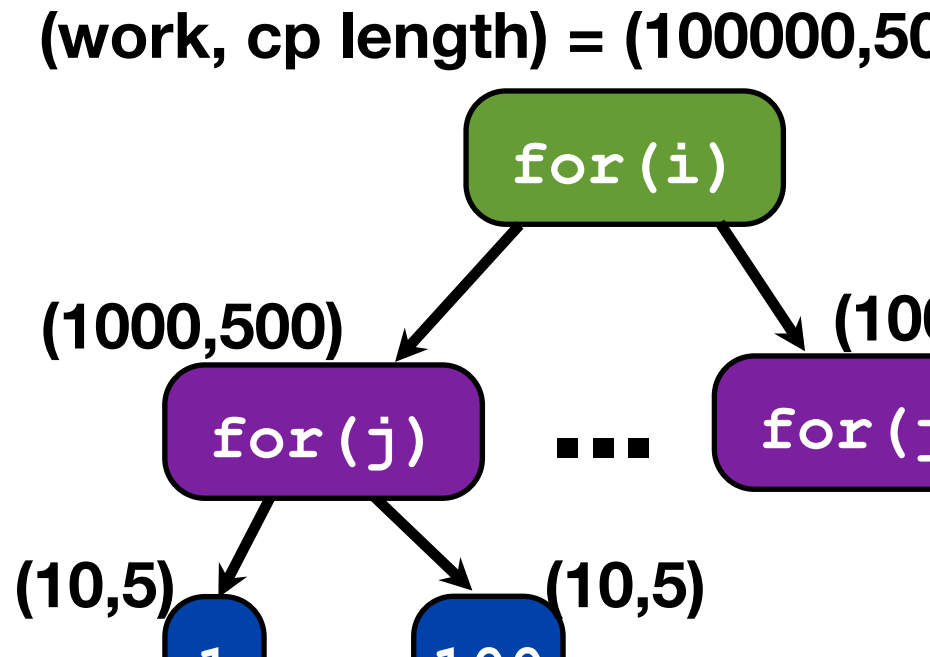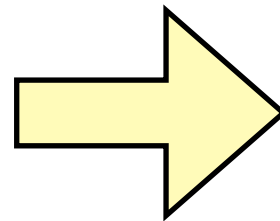
sp=2    sp=1    sp=100

**al**: Introduce localization through *region*-based analysis

adow-memory based implementation

Performs CPA analysis on every program region

Single pass: Concurrently analyzes multiple nested regions

```
1..100) {
j=1..100) {
i][j] = a[i][j-1]+3;
i][j] = b[i][j-1]+5;
```

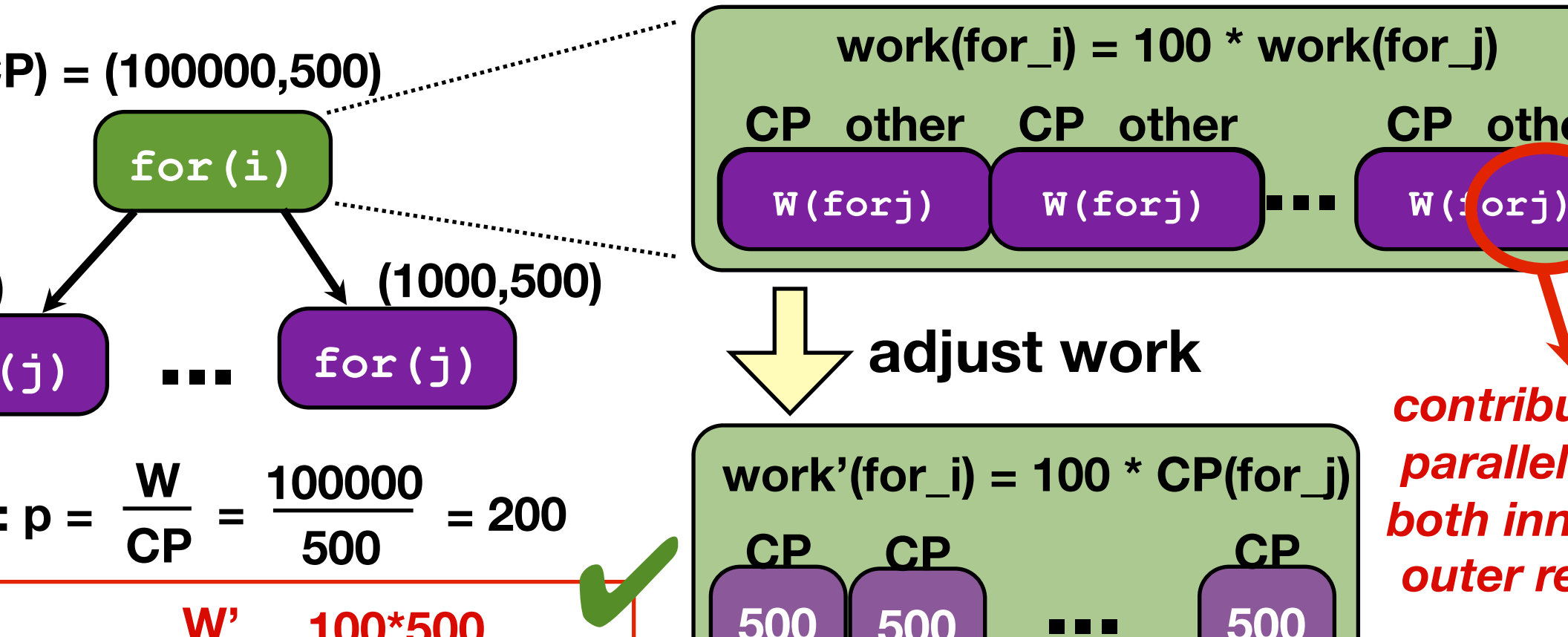**(work, cp length) = (100000,50**

**for(i)**

**(1000,500)**          **(100**

**for(j)**   **...**   **for(j**

**(10,5)**          **(10,5)**

$$p = \frac{W}{CP} = \frac{100000}{500} = 200 \; \textcolor{red}{✗}$$

**al:** Eliminate effect of nested parallelism in parallelism calculat

proximate self-parallelism using only HCPA output

"Subtracts" nested parallelism from overall parallelism

**P) = (100000,500)**

**for(i)**

**(1000,500)**

**(j)** **...** **for(j)**

$$p = \frac{W}{CP} = \frac{100000}{500} = 200$$

✓

**W'    100*500**

**work(for_i) = 100 * work(for_j)**

| CP  other | CP  other | CP  oth |
|---|---|---|
| **W(forj)** | **W(forj)** ... | **W(forj)** |

⬇ **adjust work**

*contribu
parallel
both inn
outer re*

**work'(for_i) = 100 * CP(for_j)**

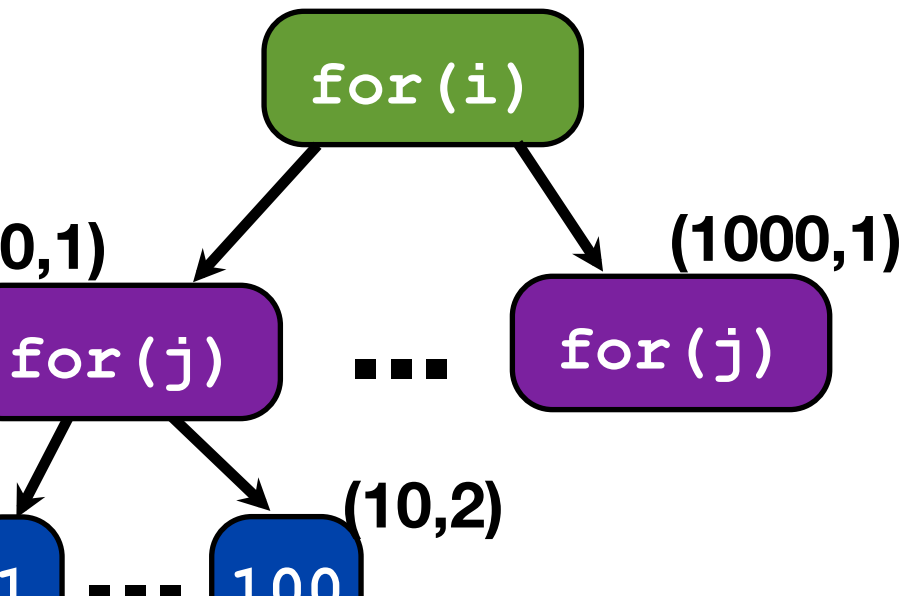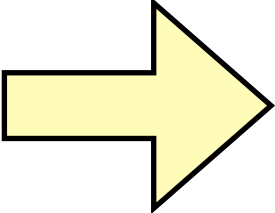| CP  CP | | CP |
|---|---|---|
| **500  500** ... | | **500** |

**al:** Convert dynamic region data to static region output

rge dynamic nodes associated with same static region

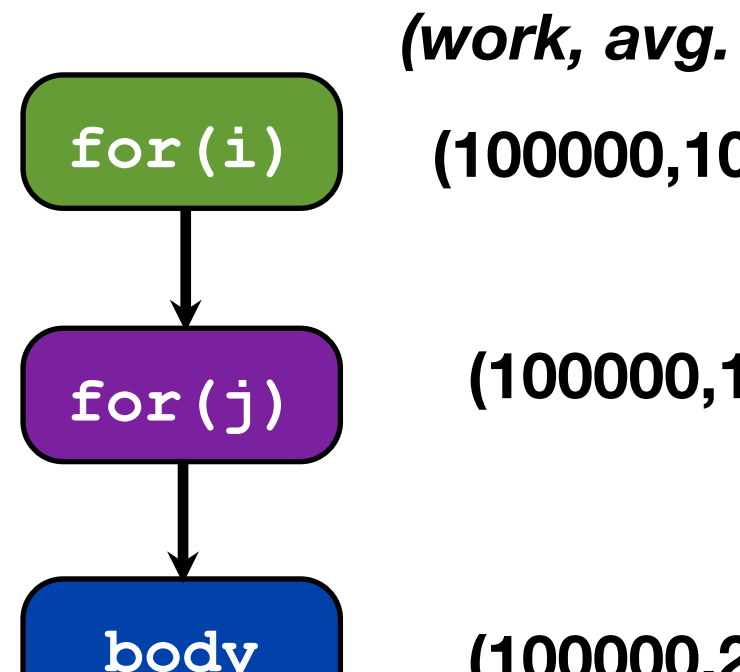Work: *Sum* of work across dynamic instances

Self-Parallelism: *Weighted Average* across dynamic instances

*(work, sp) = (100000,100)*

**for(i)**

**0,1)**

**for(j)**

**...**

**(1000,1)**

**for(j)**

**(10,2)**

**1**

**...**

**100**

**merge dynamic**

*(work, avg.*

**for(i)**     **(100000,10**

**for(j)**     **(100000,1**

**body**     **(100000,2**

# ther Details on Discovery in Our Paper

mlin handles much more complex structures than just neste
loops: finds parallelism in arbitrary code including recursion

f-parallelism metric is defined and discussed in detail in the
per

mpression technique used to reduce size of HCPA output

# ating a Parallelization Plan

**al**: Use HCPA output to select best regions for target syster

*nning personalities* allow user to incorporate system constra

Software constraints: What types of parallelism can I specify?

Hardware constraints: Synchronization overhead, etc.
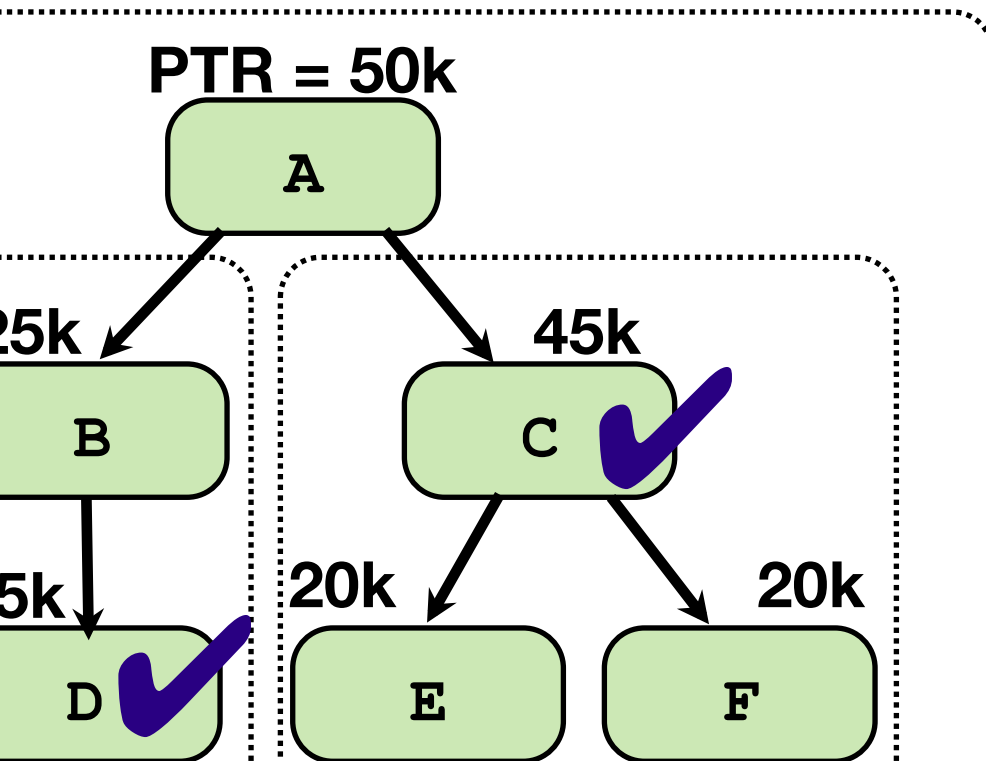
Planning algorithm can change based on constraints

# OpenMP Planner

sed on OpenMP 2.0 specification

Focused on loop-level parallelism

Disallows nested parallelism because of overhead

Planning algorithm based on dynamic programming

**PTR = 50k**

**parallelized time reduction = W - (W/**



| Region | Work | SP |
|--------|------|-----|
| A | 100k | 2 |
| B | 50k | 2 |
| C | 50k | 10 |
| D | 50k | 10 |
| E | 25k | 5 |

# luation

Ran Kremlin on serial versions; targeting OpenMP

Parallelized according to Kremlin's plan

Gathered performance results on 8 socket AMD 8380 Quad-core

Compared against third-party parallelized versions (3rd Party)

nchmarks: NAS OpenMP and SpecOMP

Have both serial and parallel versions

Wide range of parallel speedup (min: 1.85x, max: 25.89x) on 32 cor
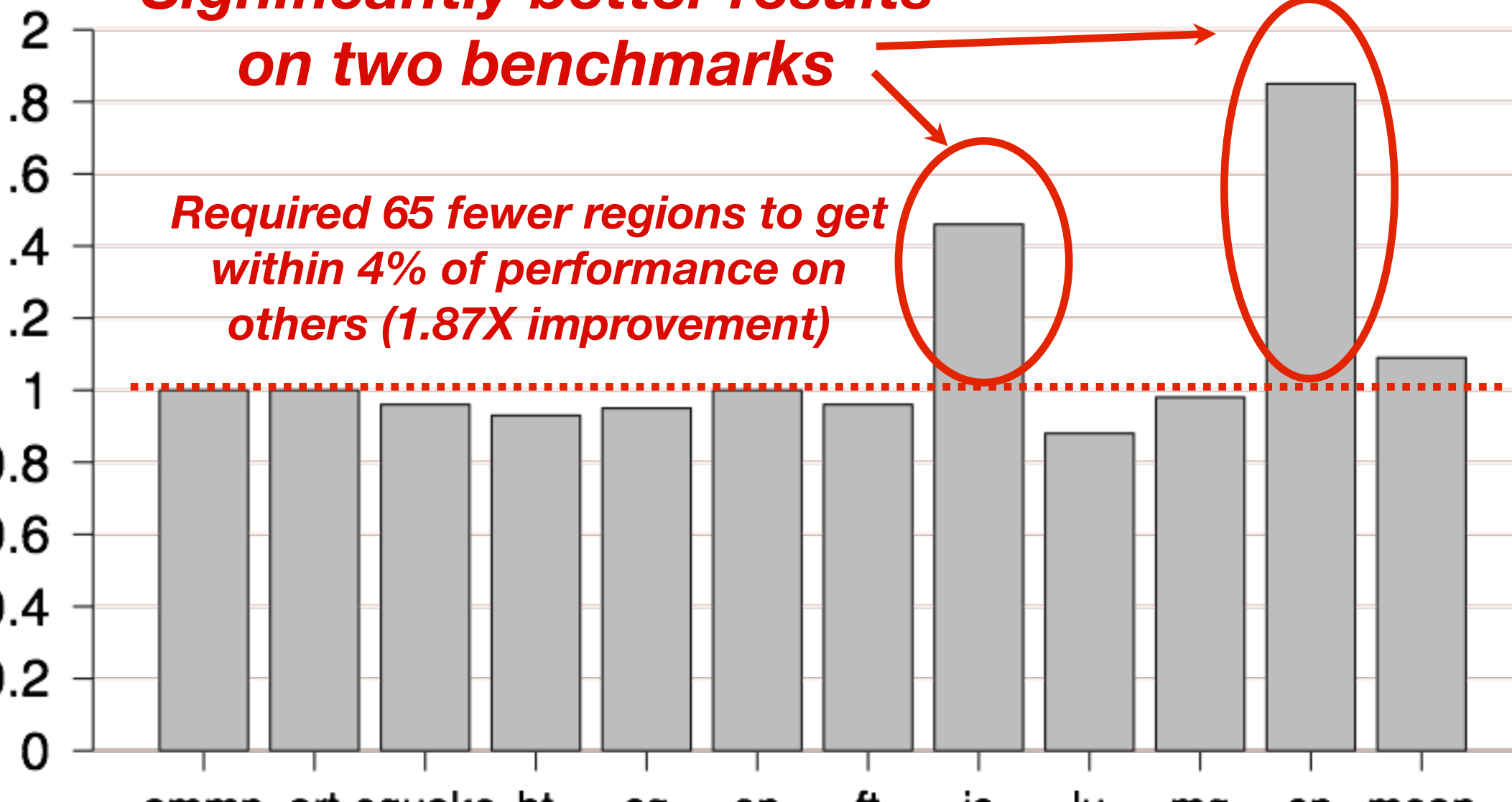
# How much effort is saved using Kremlin?

| Suite | Benchmark | # of Regions Parallelized | | Reduction |
|-------|-----------|--------------------------|---------|-----------|
| | | 3rd Party | Kremlin | |
| | art | 3 | 4 | 0.75x |
| SpecOMP | ampp | 6 | 3 | 2.00x |
| | equake | 10 | 6 | 1.67x |
| | ep | 1 | 1 | 1.00x |
| | is | 1 | 1 | 1.00x |
| | ft | 6 | 6 | 1.00x |
| NPB | mg | 10 | 8 | 1.25x |
| | cg | 22 | 9 | 2.44x |
| | lu | 28 | 11 | 2.55x |
| | bt | 54 | 27 | 2.00x |
| | sp | 70 | 58 | 1.21x |
| | **Overall** | **211** | **134** | **1.57x** |

# good is Kremlin-guided performance?

mpared performance against expert, third-party version



**Significantly better results on two benchmarks**

**Required 65 fewer regions to get within 4% of performance on others (1.87X improvement)**

ermined what % of speedup comes from first {25,50,75,100
ecommended regions

| | Fraction of Kremlin Plan Applied | | | |
|---|---|---|---|---|
| | **First 25%** of regions | **Second 25%** of regions | **Third 25%** of regions | **Last 25%** regions |
| Marginal benefit (% ax speedup) (avg) | 56.2% | 30.2% | 9.2% | 4.4% |

# onclusion

- remlin helps a programmer determine:
  - *What parts of this program should I spend time parallelizing*

- ee key techniques introduced by Kremlin
  - *Hierarchical CPA*: How much total parallelism is in each region?
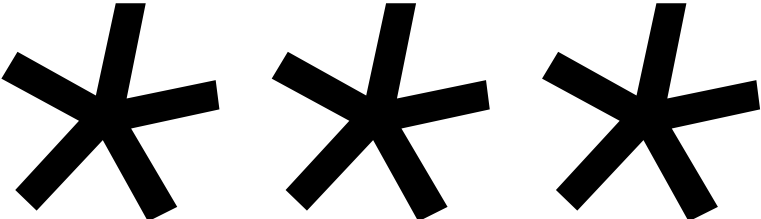  - *Self-Parallelism*: How much parallelism is only in this region?
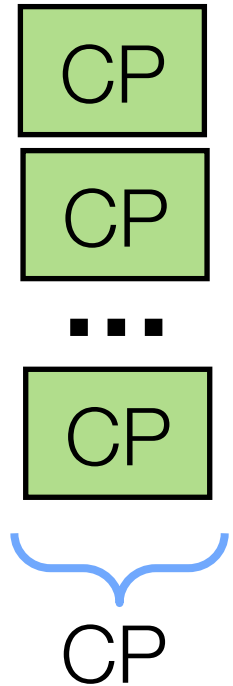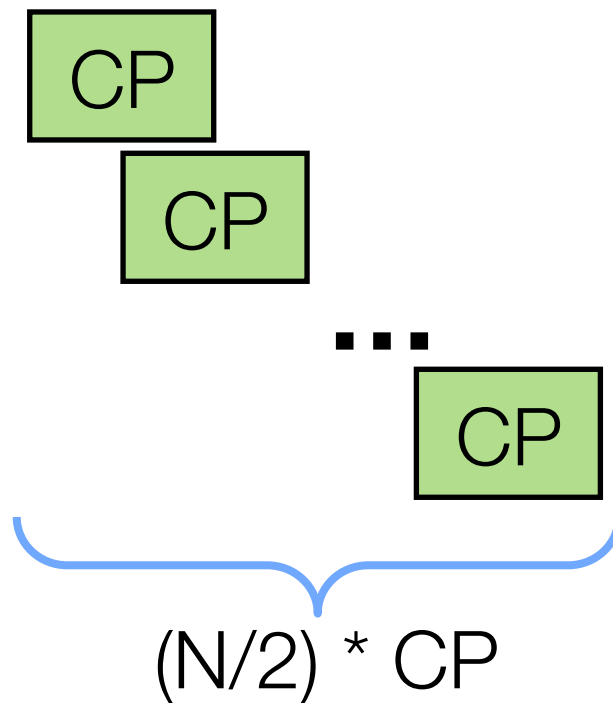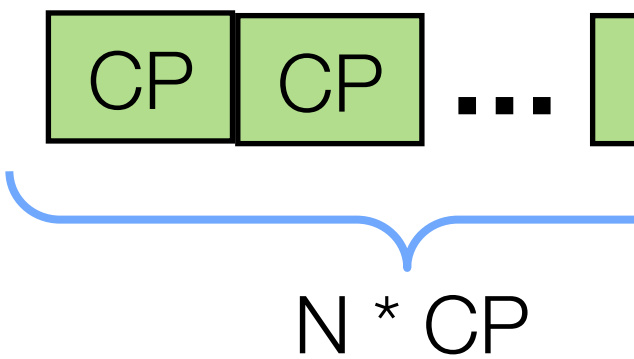  - *Planning Personalities*: What regions are best for my target system?

- mpelling results
  - .57x average reduction in number of regions parallelized
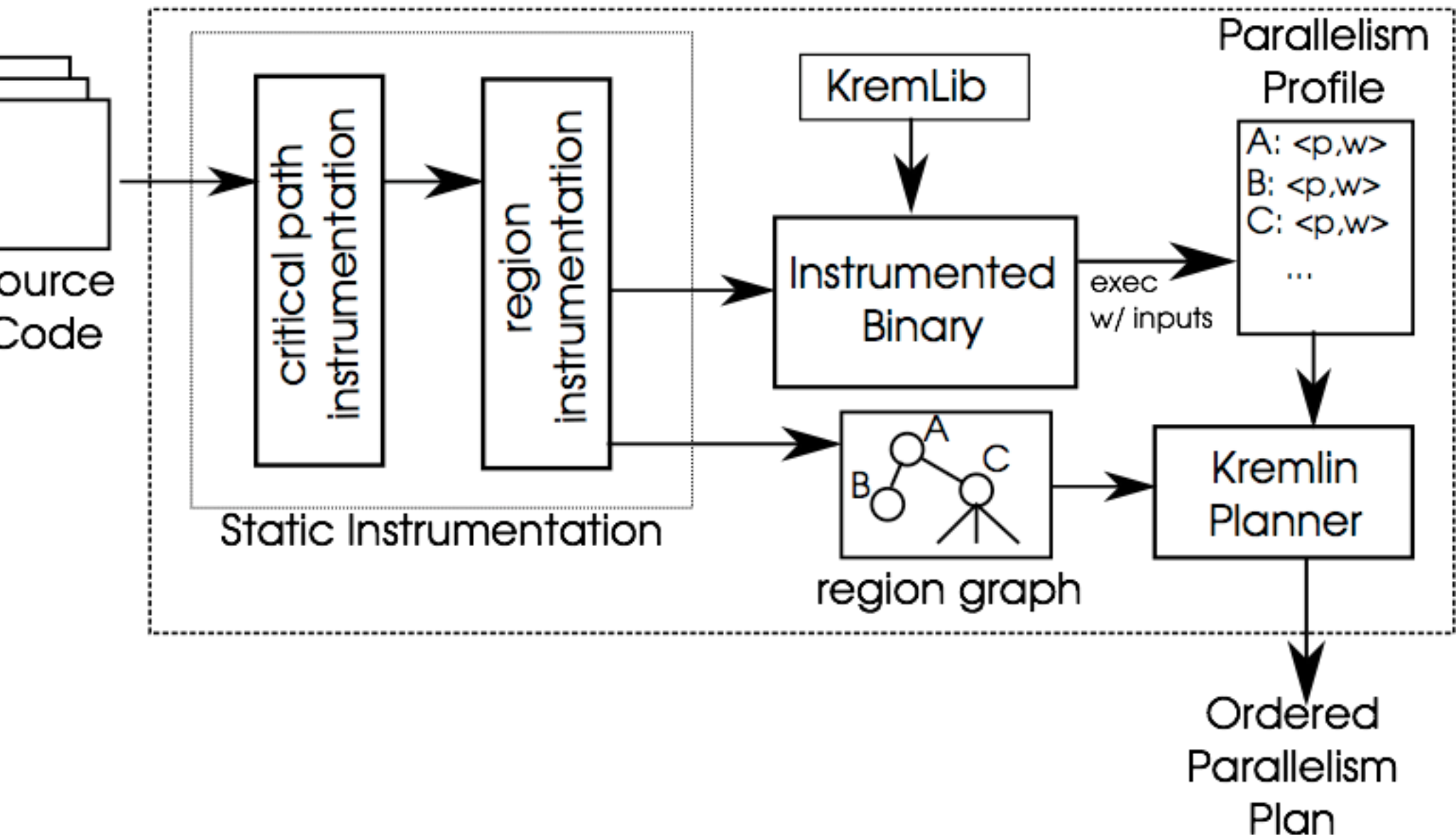  - Greatly improved performance on 2 of 11 benchmarks; very close others
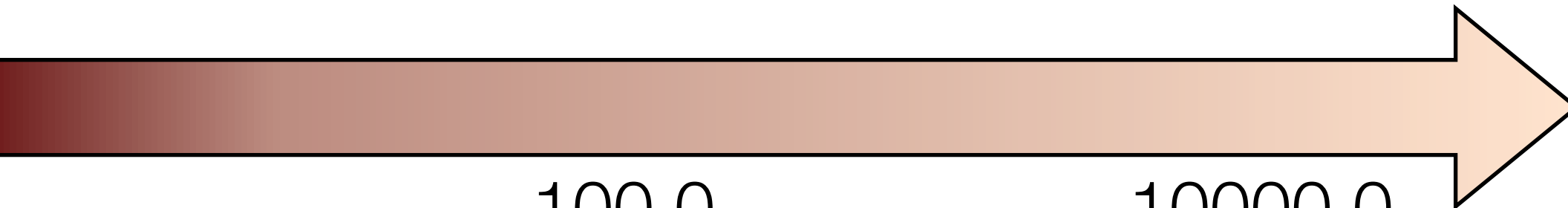
* * *

| ype | **DOALL** | **DOACROSS** | **Serial** |
|---|---|---|---|
| ath CP) | CP<br><br>CP<br><br>...<br><br>CP<br><br>CP | CP<br><br>CP<br><br>...<br><br>CP<br><br>(N/2) * CP | CP CP ... <br><br>N * CP |
| ") | N * CP | N * CP | N * CP |
| m | $\dfrac{N * CP}{CP} =$ **N** | $\dfrac{N * CP}{(N/2) * CP} =$ **2.0** | $\dfrac{N * CP}{N * CP} =$ **1.** |

Source Code → Static Instrumentation [ critical path instrumentation → region instrumentation ] → Instrumented Binary (KremLib →) —exec w/ inputs→ Parallelism Profile (A: <p,w>, B: <p,w>, C: <p,w>, ...) → Kremlin Planner → Ordered Parallelism Plan; region instrumentation → region graph (A, B, C) → Kremlin Planner

100.0                    10000.0

**ally Serial**                    **Highly Parallel**

*All* work is
*on* critical
path
(ET == CP)

*Most* w
*off* cr
pa
(ET >

***Parallelism is a result of execution time***