



# UNDERSUBSCRIBED THREADING ON CLUSTERED CACHE ARCHITECTURES

WIM HEIRMAN<sup>1,2</sup> TREVOR E. CARLSON<sup>1</sup> KENZO VAN CRAEYNES<sup>1</sup>  
IBRAHIM HUR<sup>2</sup> AAMER JALEEL<sup>2</sup> LIEVEN EECKHOUT<sup>1</sup>

<sup>1</sup> GHENT UNIVERSITY

<sup>2</sup> INTEL CORPORATION



HPCA 2014, ORLANDO, FL

# CONTEXT

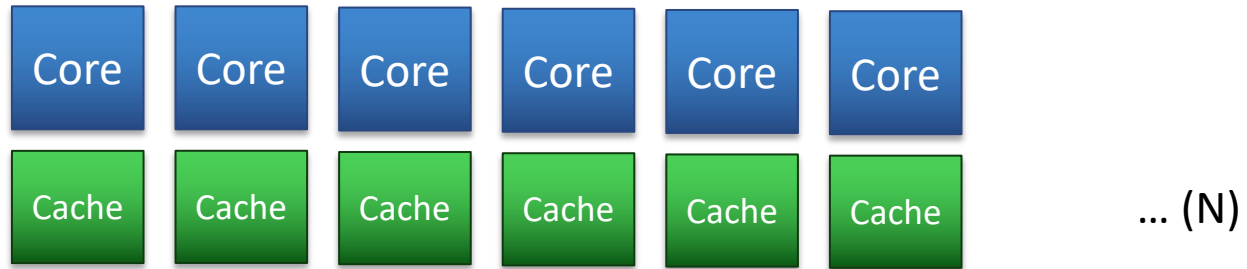
- Many-core processor with 10s-100s of cores
  - E.g. Intel Xeon Phi, Tiler, GPGPU
- Running scalable, data-parallel workloads
  - SPEC OMP, NAS Parallel Benchmarks, ...
- Processor design @ fixed area/power budget
  - Spend on cores, caches?
  - Cache topology?

# OVERVIEW

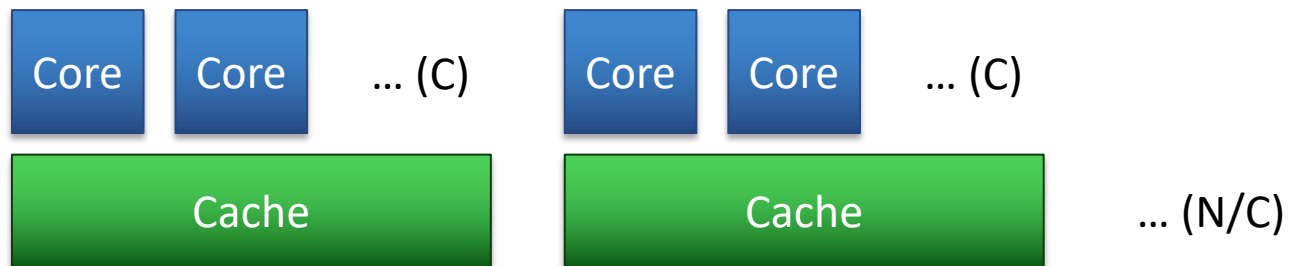
- Cache topology:
  - Why clustered caches?
  - Why undersubscription?
- Dynamic undersubscription
  - CRUST algorithms for automatic adaptation
- CRUST and future many-core design

# MANY-CORE CACHE ARCHITECTURES

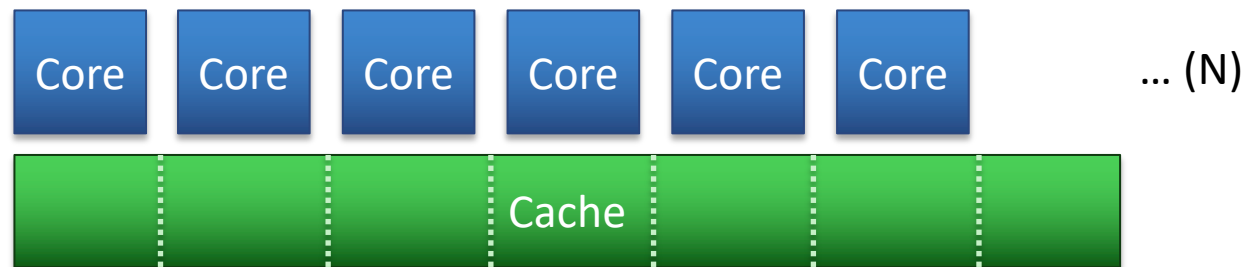
private



clustered



shared  
(NUCA)



# MANY-CORE CACHE ARCHITECTURES

private



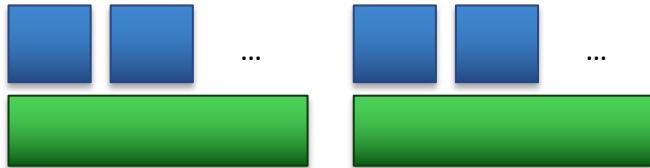
hit latency



sharing



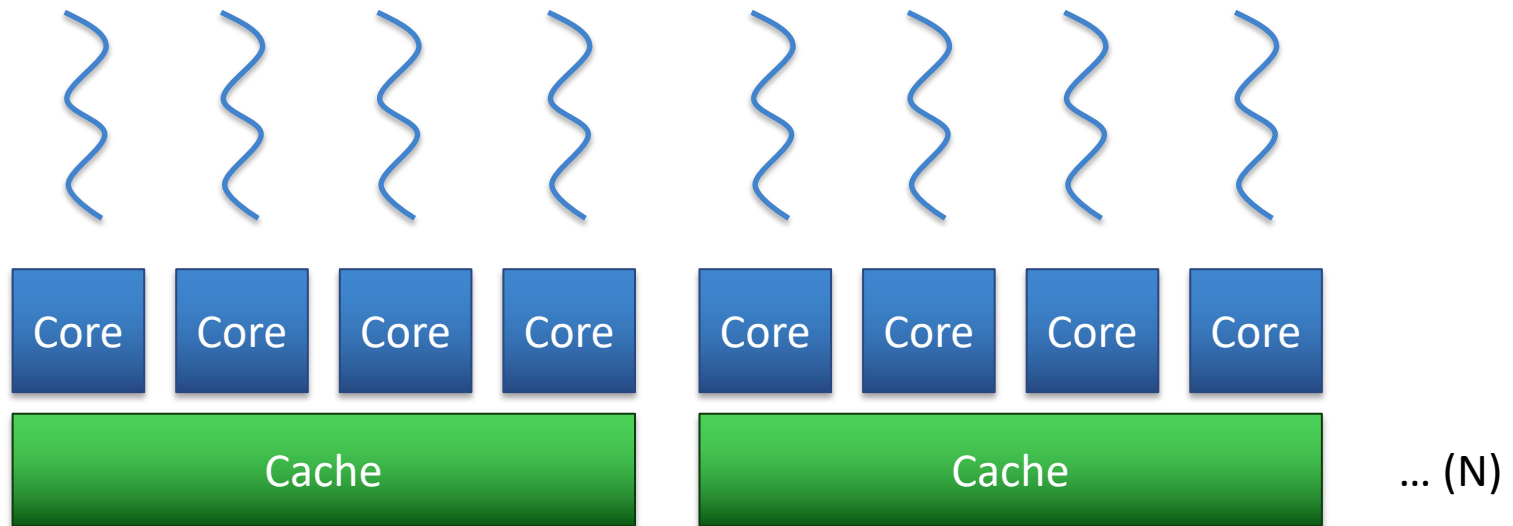
clustered



shared

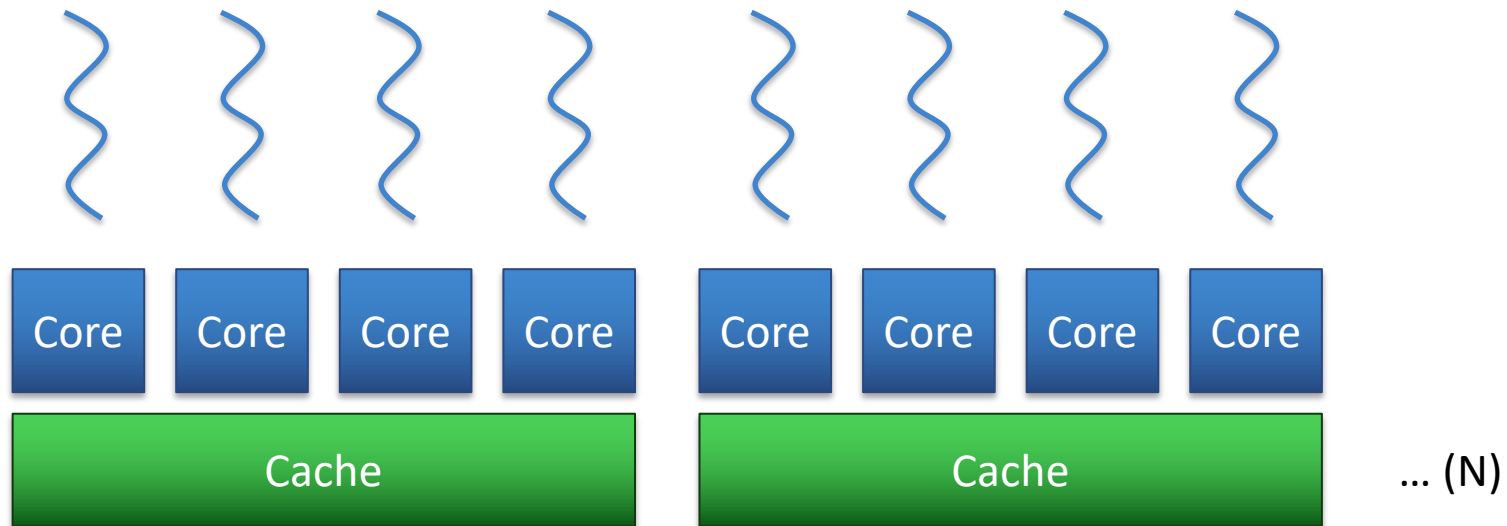


# UNDERSUBSCRIBING FOR CACHE CAPACITY



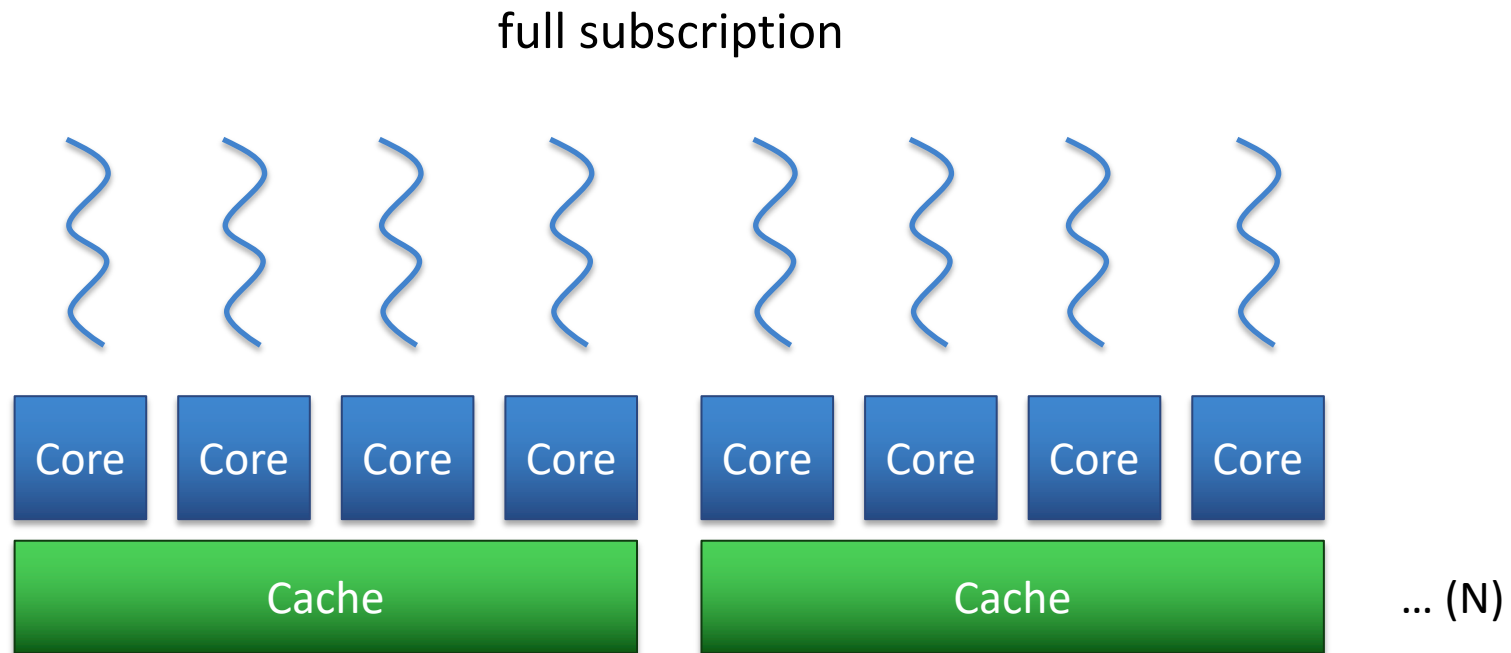
# UNDERSUBSCRIBING FOR CACHE CAPACITY

- Less than  $C$  active cores/threads per cluster



# UNDERSUBSCRIBING FOR CACHE CAPACITY

- Less than  $C$  active cores/threads per cluster

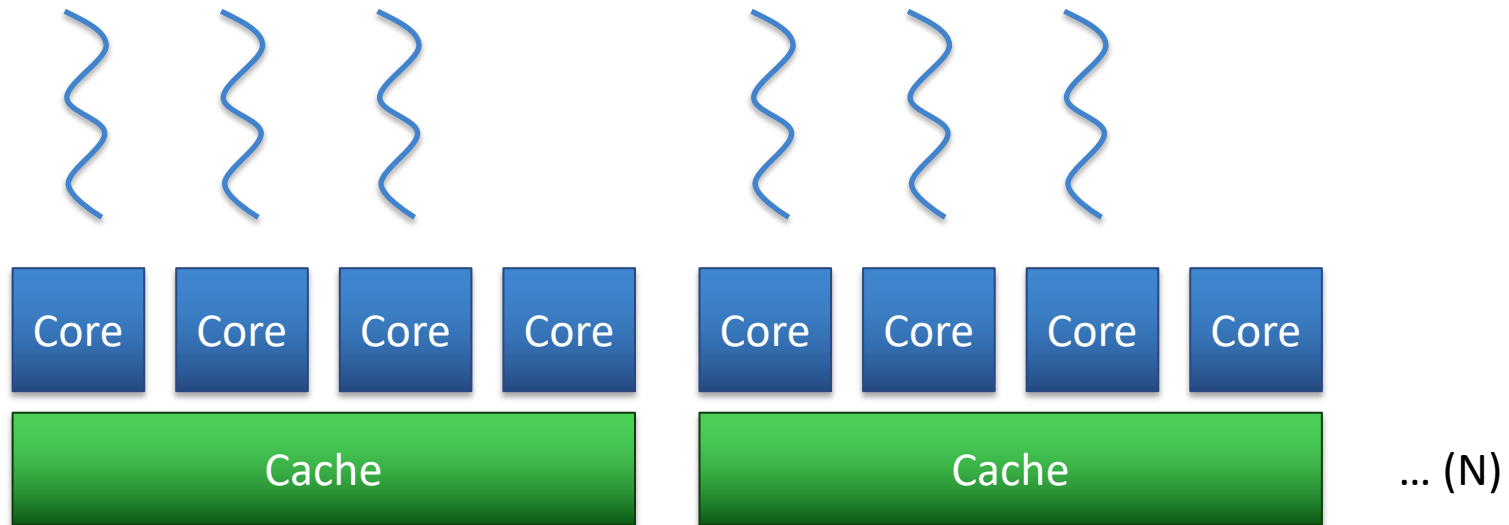




# UNDERSUBSCRIBING FOR CACHE CAPACITY

- Less than C active cores/threads per cluster

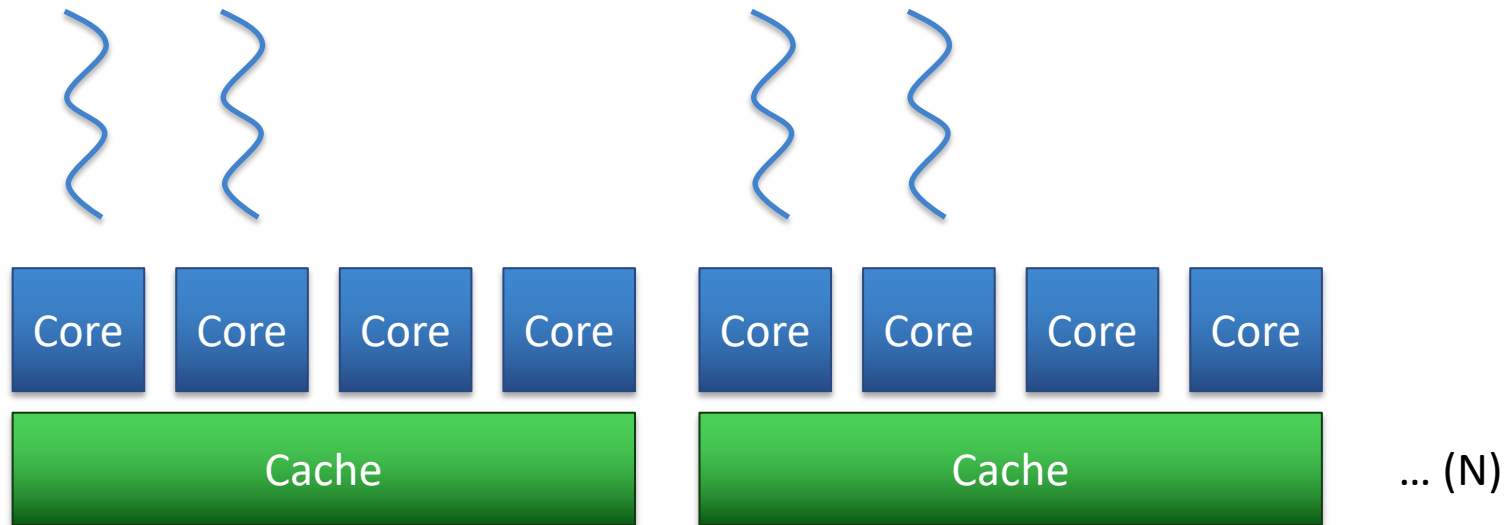
3/4 undersubscription



# UNDERSUBSCRIBING FOR CACHE CAPACITY

- Less than C active cores/threads per cluster

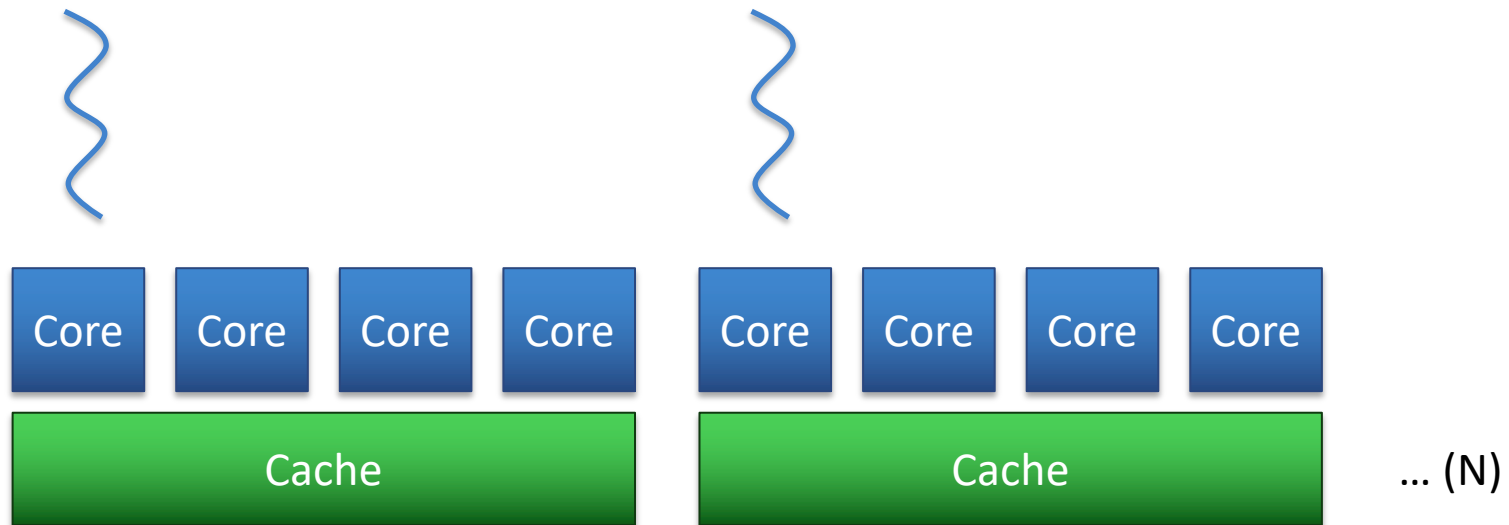
2/4 undersubscription



# UNDERSUBSCRIBING FOR CACHE CAPACITY

- Less than  $C$  active cores/threads per cluster

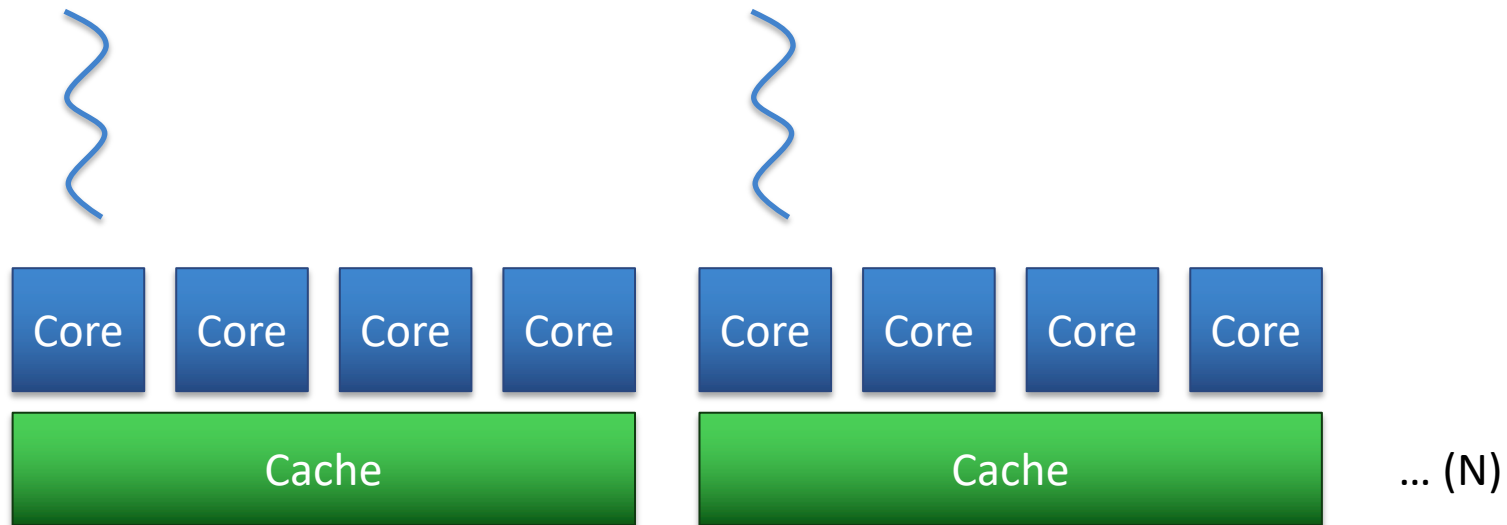
1/4 undersubscription



# UNDERSUBSCRIBING FOR CACHE CAPACITY

- Less than  $C$  active cores/threads per cluster
- When working set does not fit in cache

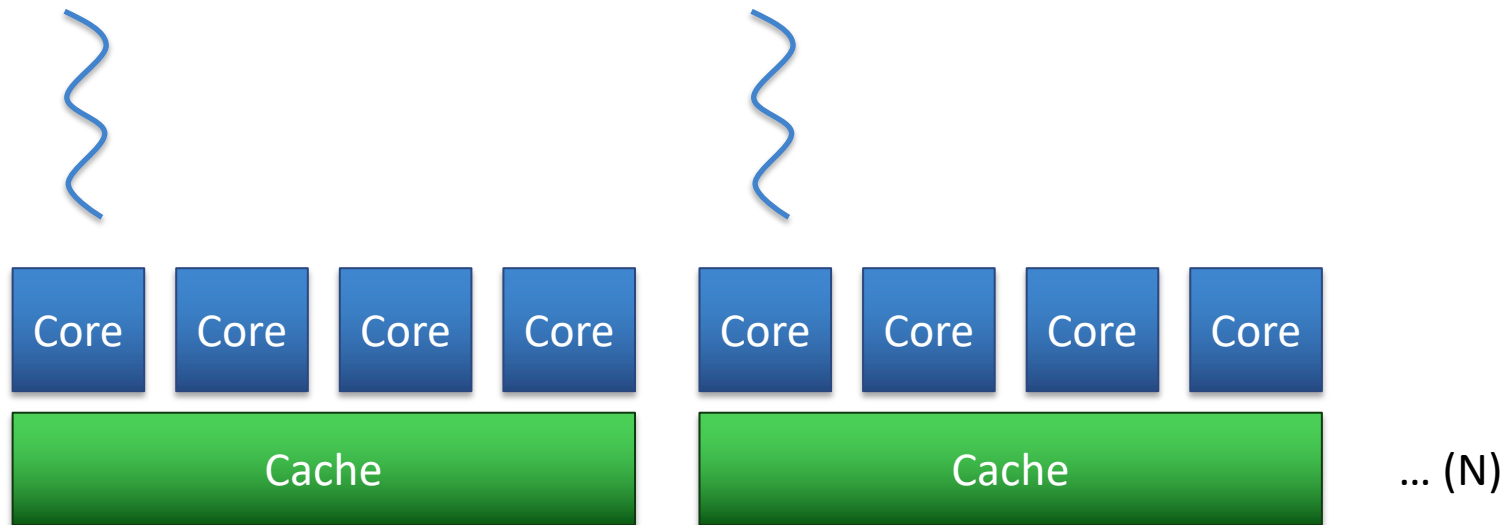
1/4 undersubscription



# UNDERSUBSCRIBING FOR CACHE CAPACITY

- Less than  $C$  active cores/threads per cluster
- When working set does not fit in cache
- *Keep all cache capacity accessible*

1/4 undersubscription



# MANY-CORE CACHE ARCHITECTURES

private



hit latency



sharing

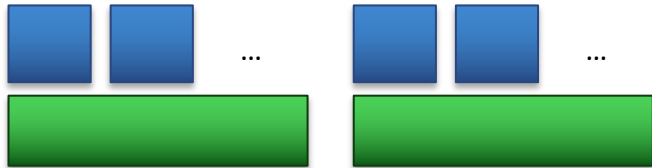


undersubscription



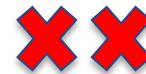
(1:1)

clustered



(1:C)

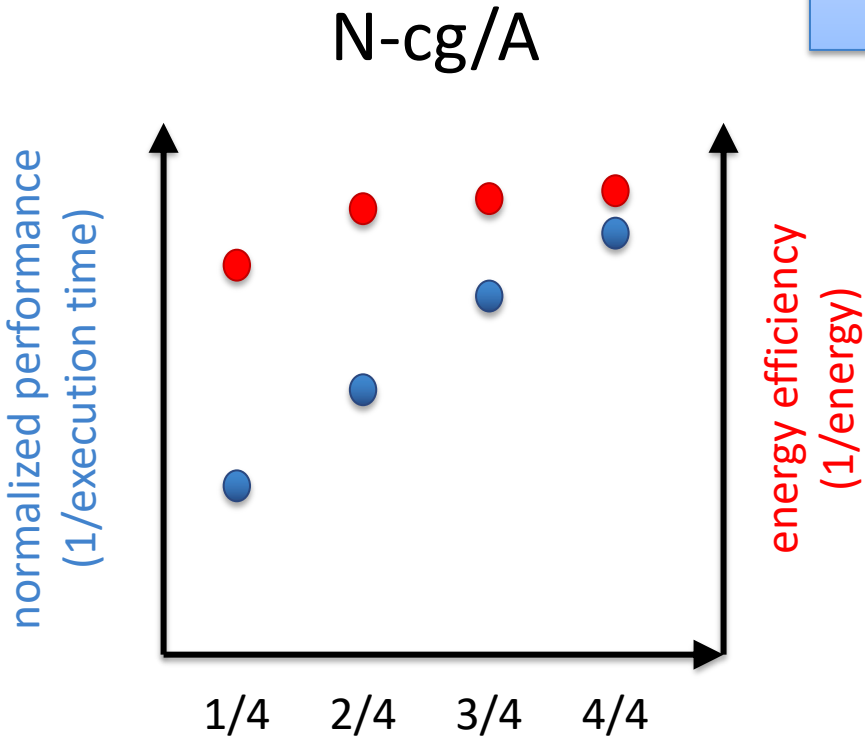
shared



(1:N)

# PERFORMANCE & ENERGY: WORKING SET VS. CACHE SIZE

- Baseline architecture:
- 128 cores
  - private L1
  - clustered L2  
1M shared per 4 cores



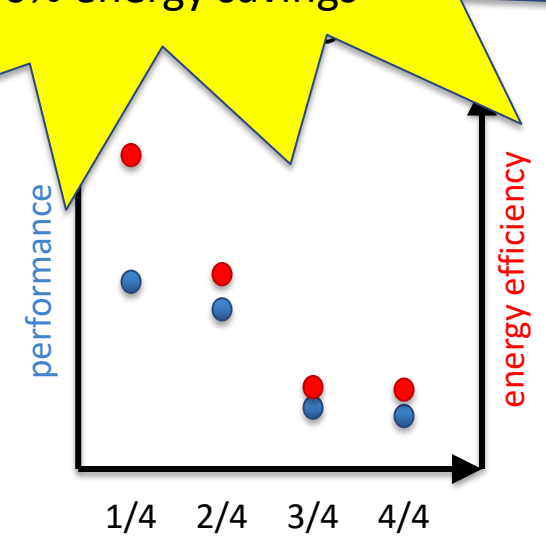
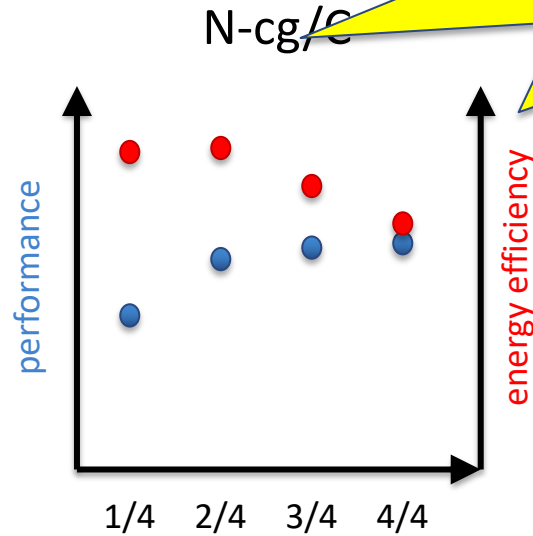
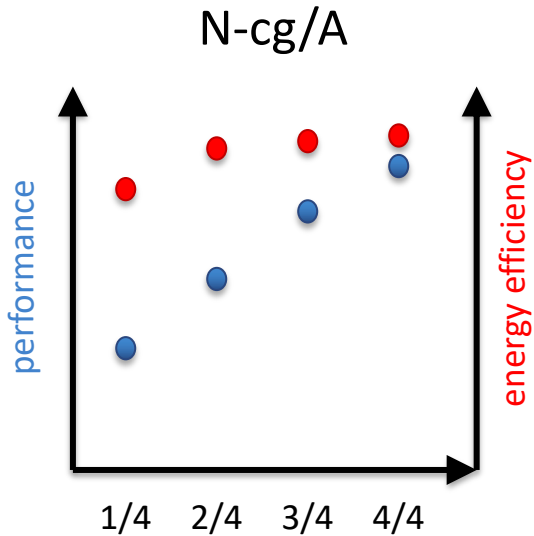
# PERFORMANCE & ENERGY: WORKING SET VS. CACHE SIZE

**Compute bound:**  
use all cores for  
highest performance

**Bandwidth bound:**  
disable cores for better  
energy efficiency

**Capacity bound:**  
reduce thread count  
to optimize hit rate

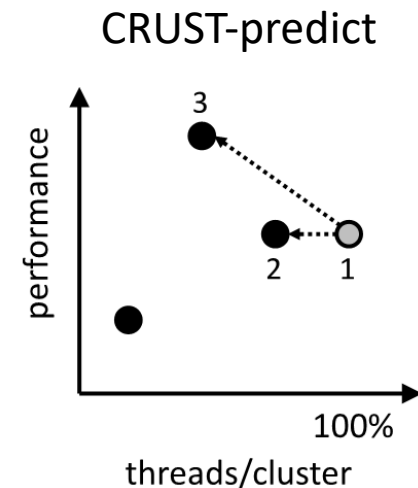
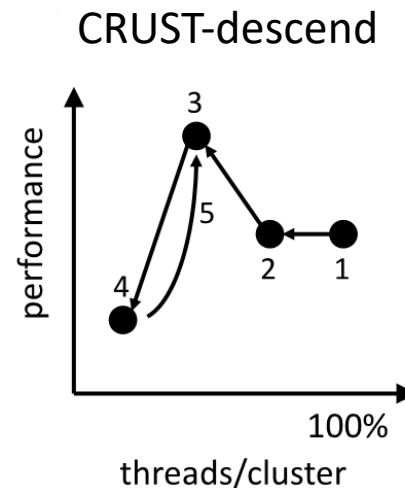
1/4 undersubscription:  
3.5x performance  
80% energy savings





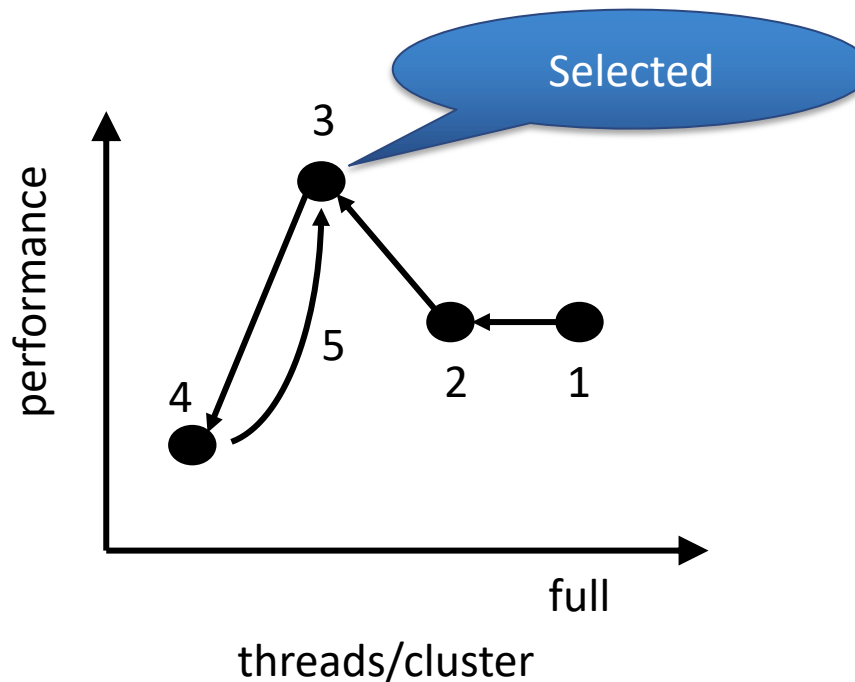
# CLUSTER-AWARE UNDERSUBSCRIBED SCHEDULING OF THREADS (CRUST)

- Dynamic undersubscription
- Integrated into the OpenMP runtime library
  - Adapt to each *#pragma omp* individually
- Optimize for performance first, save energy when possible
  - Compute bound: full subscription
  - Bandwidth bound: no\* performance degradation (\* <5% vs. full)
  - Capacity bound: highest performance
- Two CRUST heuristics (*descend* and *predict*) for on-line adaptation



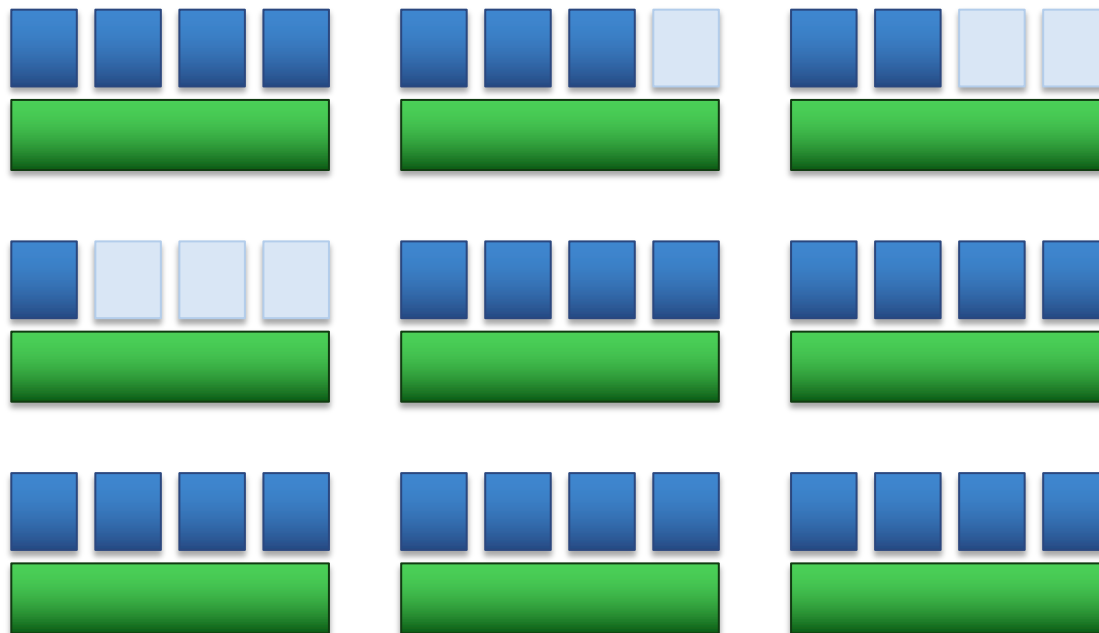
# CRUST-DESCEND

- Start with full subscription
- Reduce thread count while performance increases



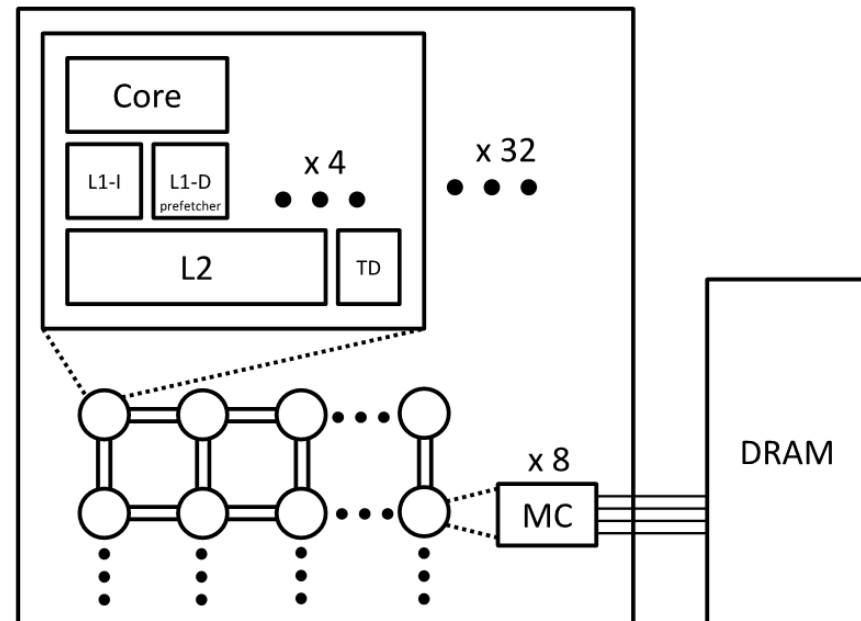
# CRUST-PREDICT

- Reduce number of steps required by being smarter
- Start with heterogeneous undersubscription
  - Measure LLC miss rate for each thread/cluster option
  - Predict performance of each option using PIE-like model
- Select best predicted option



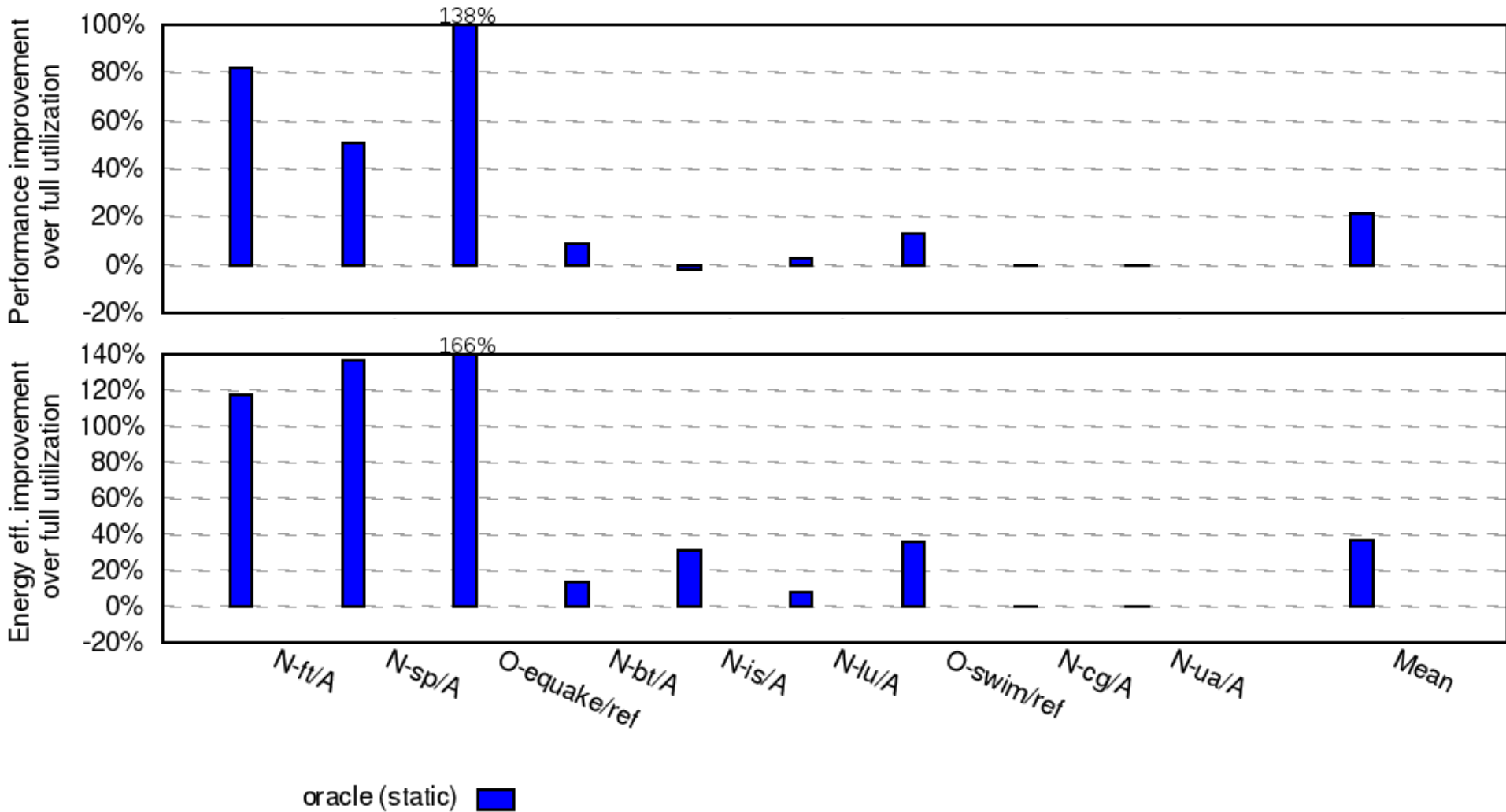
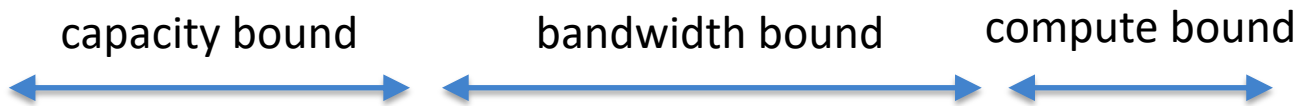
# METHODOLOGY

- Generic many-core architecture
  - 128 cores, 2-issue OOO @1GHz
  - 2x 32 KB private L1 I+D
  - L1-D stride prefetcher
  - **1 MB shared L2 per 4 cores**
  - 2-D mesh NoC
  - 64 GB/s total DRAM bandwidth

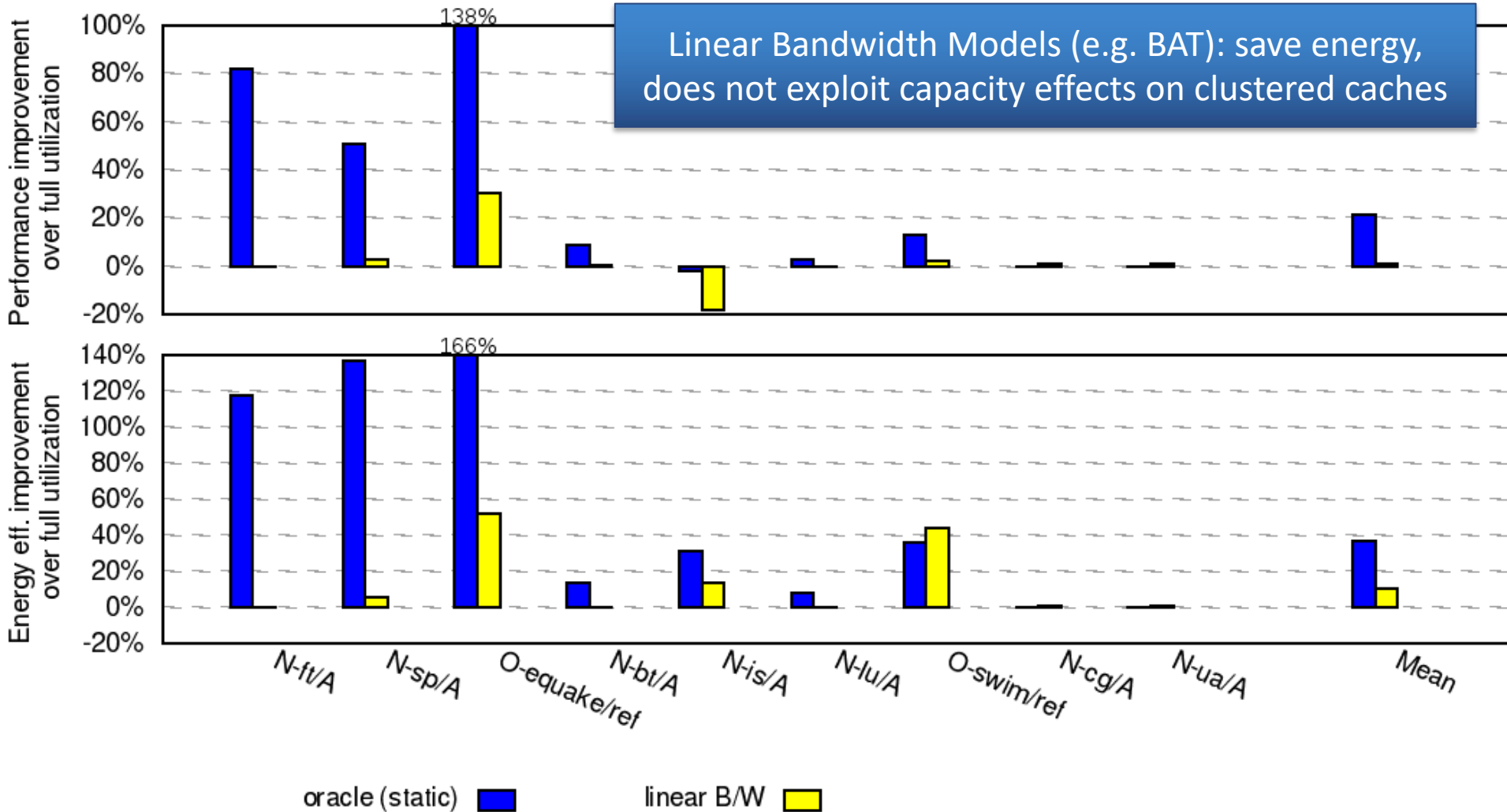
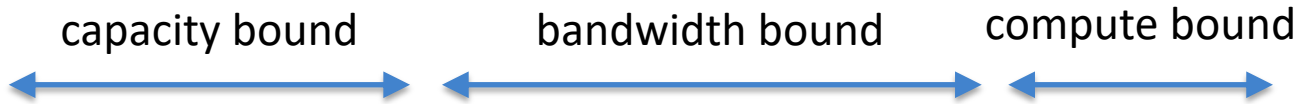


- Sniper simulator, McPAT for power
- SPEC OMP and NAS parallel benchmarks
  - Reduced iteration counts from ref, class A inputs

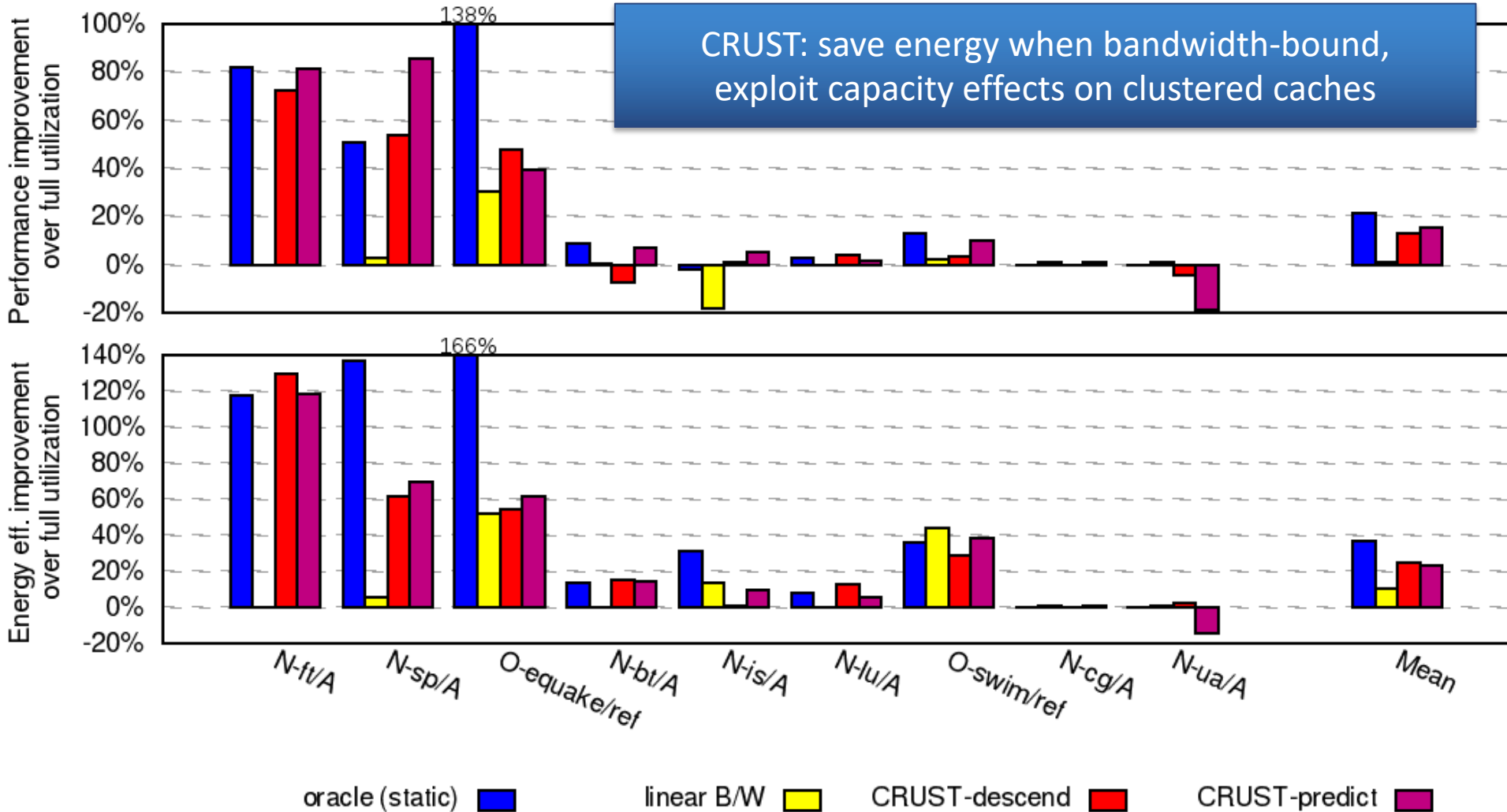
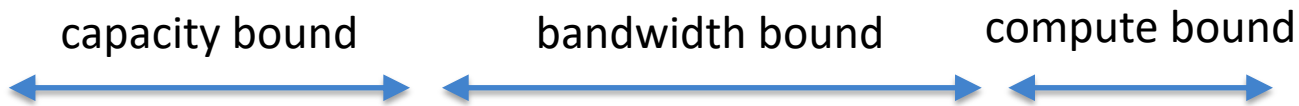
# RESULTS: ORACLE (STATIC)



# RESULTS: LINEAR BANDWIDTH MODELS



# RESULTS: CRUST



# UNDERSUBSCRIPTION VS. FUTURE DESIGNS

- Finite chip area, spent on cores or caches
  - *Increasing max. compute vs. keeping cores fed with data*
- Undersubscription can adapt workload behavior to the architecture

*Does this allow us to build a higher-performance design?*

- Sweep core vs. cache area ratio for 14-nm design
  - Fixed 600 mm<sup>2</sup> area, core = 1.5 mm<sup>2</sup>, L2 cache = 3 mm<sup>2</sup>/MB
  - Clustered L2 shared by 4 cores, latency ~ log<sub>2</sub>(size)
  - 1 GB @ 512 GB/s on-package, 64 GB/s off-package

Variant	A	B	C	D	E	F
Cores	96	128	160	192	224	256
L2 size (MB/core)	1.5	1.0	0.8	0.5	0.4	0.3
Core area	25%	33%	40%	50%	58%	64%

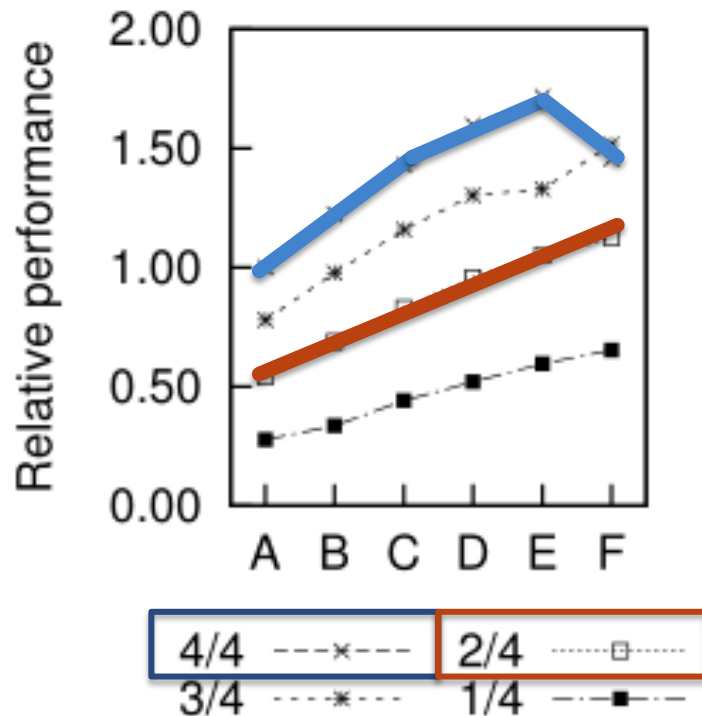


# UNDERSUBSCRIPTION FOR FUTURE DESIGNS

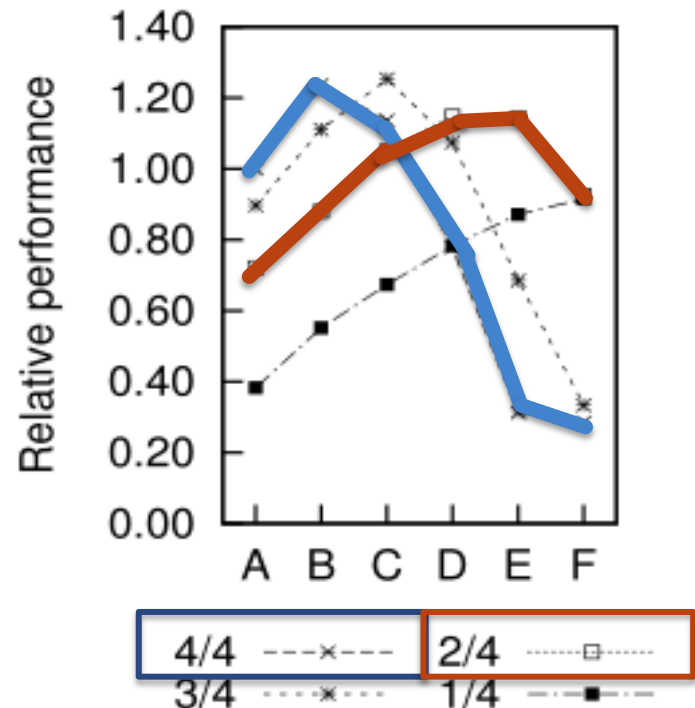
**Compute bound:** linear relation between active cores and performance

**Capacity bound:** reduce thread count until combined working set fits available cache

O-ampp/ref

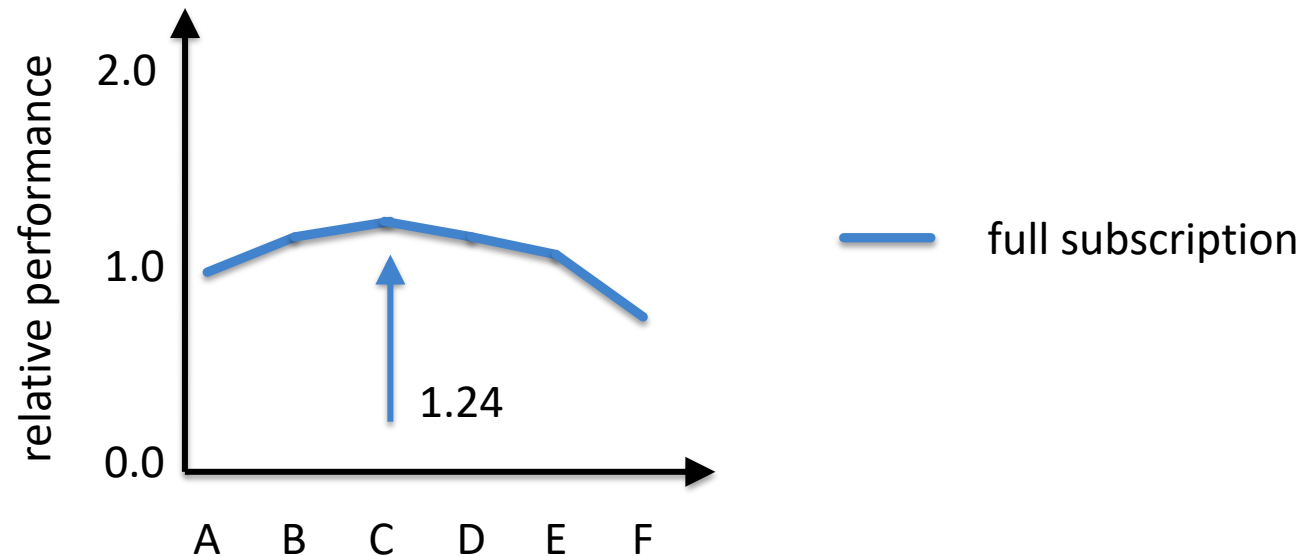


N-ft/C



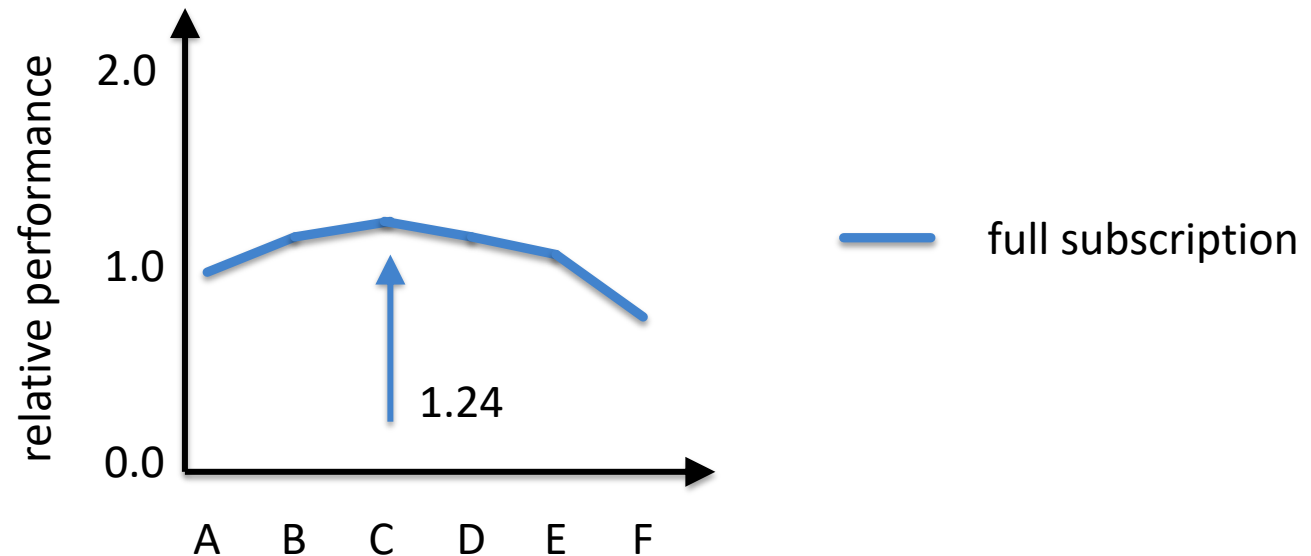
# UNDERSUBSCRIPTION FOR FUTURE DESIGNS

- Build *one* design with best average performance



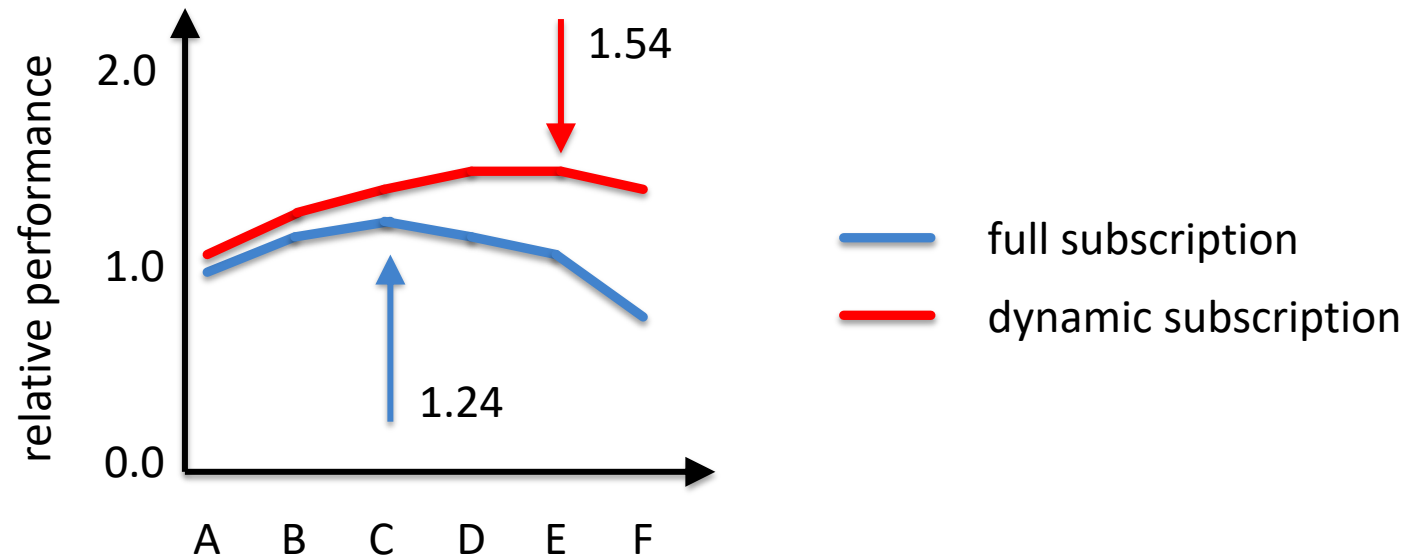
# UNDERSUBSCRIPTION FOR FUTURE DESIGNS

- Build *one* design with best average performance
- Full subscription:
  - conservative option C has highest average performance



# UNDERSUBSCRIPTION FOR FUTURE DESIGNS

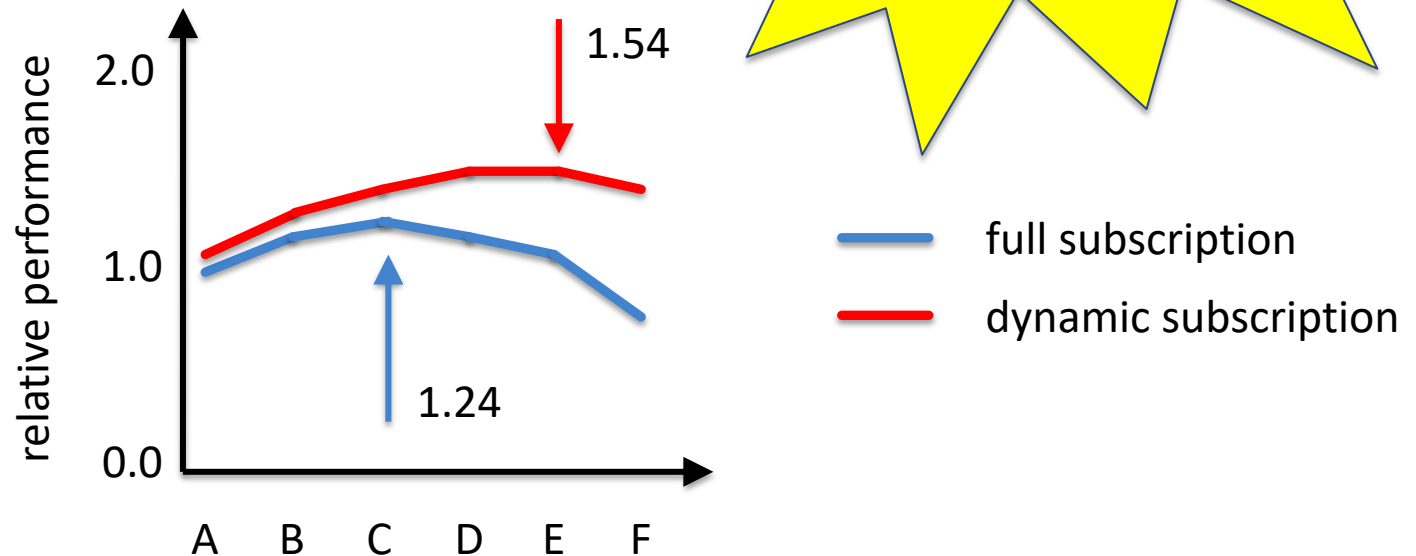
- Build *one* design with best average performance
- Full subscription:
  - conservative option C has highest average performance
- Dynamic undersubscription: prefer more cores
  - higher max. performance for compute-bound benchmarks
  - use undersubscription to accommodate capacity-bound workloads



# UNDERSUBSCRIPTION FOR FUTURE DESIGNS

- Build *one* design with best average performance
- Full subscription:
  - conservative option C has highest average performance
- Dynamic undersubscription: prefer E
  - higher max. performance for compute
  - use undersubscription to accommodate *code*

E vs. C: 40% more cores,  
15% higher performance



# CONCLUSIONS

- Use clustered caches for future many-core designs
  - Balance hit rate and hit latency
  - Exploit sharing to avoid duplication
  - Allow for undersubscription (use all cache, not all cores)
- CRUST for dynamic undersubscription
  - Adapt thread count per OpenMP parallel section
  - Performance and energy improvements of up to 50%
- Take undersubscription usage model into account when designing future many-core processors
  - CRUST-aware design: 40% more cores, 15% higher performance