# Engineering Robust Server Software Scalability





### • What does scalability mean?







### Intro To Scalability

#### • What does scalability mean?



#### Low scale-ability















- What does **scalability** mean?
  - How does performance change with resources?









- What does scalability mean?
  - How does performance change with resources?
  - How does performance change with load?





# **Scalability Terms**

- Scale Out: Add more nodes
  - More computers
- Scale Up: Add more stuff in each node
  - More processors in one node
- **Strong Scaling**: How does time change for fixed problem size?
  - Do 100M requests, add more cores -> speedup?
- Weak Scaling: How does time change for fixed (problem size/core)?
  - Do (100\*N)M requests, with N cores -> speedup?











# Speedup (N) = $\frac{S + P}{S + \frac{P}{N}}$



### Amdahl's Law









### Amdahl's Law



#### Parallel Portion: P = 4





• 10/8 = 1.25x speedup = 25% increase in **throughput**. 8/10 = 0.8x = 20% reduction in **latency** hike

### **Amdahl's Law**













• 10/7 = 1.42x speedup = 42% increase in throughput. 7/10 = 0.7x = 30% reduction in **latency** nike

### **Amdahl's Law**











• 10/6 = 1.67x speedup = 67% increase in throughput. 6/10 = 0.6x = 40% reduction in **latency** Juke

### Amdahl's Law











 Anne Bracy: "Don't try to speed up brushing your teeth" • What does she mean?



### Amdahl's Law





- Why don't we get (Nx) speedup with N cores?
  - What prevents ideal speedups?





# Impediments to Scalability

- Shared Hardware
  - Functional Units
  - Caches
  - Memory Bandwidth
  - IO Bandwidth
  - . . .
- Data Movement
  - From one core to another
- Blocking
  - Locks (and other synchronization)
  - Blocking IO



# Impediments to Scalability

- Shared Hardware
  - Functional Units
  - Caches
  - Memory Bandwidth
  - IO Bandwidth

- Data Movement
  - From one core to another
- Blocking
  - Blocking IO
- Juke Locks (and other synchronization)

#### Let's talk about these for now











#### A core has 2 threads (2-way SMT) - Also private L1 + L2 caches (not shown)











#### 4 cores share an LLC

- Connected by on chip interconnect



# **Hypothetical System**





#### We have a 2 socket node

- Has 2 chips
- DRAM
- Also some IO devices (not shown)





#### We have 2 nodes







Suppose we have 2 requests: where best to run them?



#### Different threads on same core?





#### Different cores on same chip?







#### Different chips on same node?









#### Different nodes?





# How To Control Placement?

- Within a node: sched\_setaffinity
  - Set mask of CPUs that a thread can run on
  - SMT contexts have different CPU identifiers
  - In pthreads, library wrapper: pthread\_setaffinity\_np
- Across nodes: depends..
  - Daemons running on each node? Direct requests to them
  - Startup/end new services? Software management
  - Load balancing becomes important here





# **Tradeoff: Contention vs Locality**

#### **Increasing Communication Latency**



#### • Trade off:

- Contend for shared resources?
- Longer/slower communication?



Same Node

**Different Node** 

**Increasing Contention** 



# **Tradeoff: Contention vs Locality**

#### Increasing Communication Latency

# Same Core

Same Chip

Loads + Stores Same Cache 1s of cycles Loads + StoresLoads + StoresIO OperationsOn Chip CoherenceOff Chip CoherenceNetwork10s of cycles100s cyclesKs-Ms of cycles



Same Node

Different Node







- Non Uniform Memory Access (NUMA—technically, ccNUMA) Memory latency differs depending on physical address
- migrate\_pages, mbind: control physical memory placement



### NUMA



# **Tradeoff: Contention vs Locality**

Same Core

Chip Same



Node Same **Different Node** 

#### Increasing Contention







# **Tradeoff: Contention vs Locality**

- On chip b/w
- LLC capacity
- On chip cooling
- L1/L2 capacity - Functional Units

Core Same

Chip Same



- External network b/w - Datacenter cooling
- Memory b/w
- Chip<-> chip b/w
- 10 b/w

Node Same

Node Different

#### **Increasing Contention**





- Suppose two threads need + are sensitive to:
  - LLC Capacity
  - Memory bandwidth
- What happens when we run them together?





- Suppose two threads need + are sensitive to:
  - LLC Capacity
  - Memory bandwidth
- What happens when we run them together?
  - Contention for LLC -> more cache misses
    - Slows down program, but also...





- Suppose two threads need + are sensitive to:
  - LLC Capacity
  - Memory bandwidth
- What happens when we run them together?
  - Contention for LLC -> more cache misses
    - Slows down program, but also...
  - Increases memory bandwidth demands
    - Which we already need and are contending for :(



![](_page_33_Picture_11.jpeg)

- Suppose two threads need + are sensitive to:
  - LLC Capacity
  - Memory bandwidth
- What happens when we run them together?
  - Contention for LLC -> more cache misses
    - Slows down program, but also...
  - Increases memory bandwidth demands
    - Which we already need and are contending for :(
- Interactions can make contention even worse!
  - Is there a flip side?

![](_page_34_Picture_11.jpeg)

![](_page_34_Picture_13.jpeg)

# Improved Utilization

- Can improve utilization of resources
  - One thread executes while another stalls
  - One thread uses FUs that the other does not need
  - Pair large cache footprint with small cache footprint
  - Shared code/data: one copy in cache

![](_page_35_Picture_6.jpeg)

![](_page_35_Picture_8.jpeg)

#### • So how do we improve things?

![](_page_36_Picture_2.jpeg)

![](_page_36_Picture_5.jpeg)

#### • So how do we improve things?

- Profile our system! Understand what is slow and why
- Remember: Ahmdal's law!

![](_page_37_Picture_4.jpeg)

![](_page_37_Picture_8.jpeg)

# Intro To Scalability

![](_page_38_Figure_1.jpeg)

**Competing Requests** 

#### This graph hides a lot of important detail

![](_page_38_Picture_4.jpeg)

![](_page_38_Picture_7.jpeg)

# Intro To Scalability

![](_page_39_Figure_1.jpeg)

**Competing Requests** 

#### Breaking it down shows WHERE to focus our optimization efforts

![](_page_39_Picture_4.jpeg)

![](_page_39_Picture_6.jpeg)

- So how do we improve things?
  - Profile our system! Understand what is slow and why
  - Remember: Ahmdal's law!
- After making a change, what do we do?

![](_page_40_Picture_5.jpeg)

![](_page_40_Picture_10.jpeg)

- So how do we improve things?
  - Profile our system! Understand what is slow and why
  - Remember: Ahmdal's law!
- After making a change, what do we do?
  - Measure impact: did we make things better? How much?

![](_page_41_Picture_6.jpeg)

![](_page_41_Picture_9.jpeg)

#### • So what can we do?

- Optimize code to improve its performance
- Transform code to improve resource usage (e.g. cache space)
- Pair threads with complementary resource usage

![](_page_42_Picture_5.jpeg)

![](_page_42_Picture_7.jpeg)

- So what can we do?
  - Optimize code to improve its performance
  - Transform code to improve resource usage (e.g. cache space)
  - Pair threads with complementary resource usage
- Sounds complicated?
  - Learn more about hardware (e.g., ECE 552)
  - Take Performance/Optimization/Parallelism

![](_page_43_Picture_8.jpeg)

![](_page_43_Picture_10.jpeg)

# Impediments to Scalability

- Shared Hardware
  - Functional Units
  - Caches
  - Memory Bandwidth
  - IO Bandwidth
- Data Movement
  - From one core to another
  - Blocking

Blocking IO

Duke • Locks (and other synchronization)

#### Let's talk about this next

![](_page_44_Picture_13.jpeg)

![](_page_45_Picture_0.jpeg)

### Critical principle: never block

• Why not?

![](_page_45_Picture_3.jpeg)

![](_page_45_Picture_5.jpeg)

- Critical principle: never block
  - Why not?
- Can't we just throw more threads at it?
  - One thread per request (or even a few per request)
  - Just block whenever you want

![](_page_46_Picture_6.jpeg)

![](_page_46_Picture_7.jpeg)

![](_page_46_Picture_9.jpeg)

### Never Block

- Critical principle: never block
  - Why not?
- Can't we just throw more threads at it?
  - One thread per request (or even a few per request)
  - Just block whenever you want
- Nice in theory, but has overheads
  - Context switching takes time
  - Switching threads reduces temporal locality
    - Threads not blocked? May thrash if too many
- Threads use resources

![](_page_47_Picture_12.jpeg)

# Non-Blocking IO

### • IO operations often block (we never want to block)

• Can use non-blocking IO

![](_page_48_Picture_3.jpeg)

![](_page_48_Picture_5.jpeg)

# **Non-Blocking IO**

- IO operations often block (we never want to block)
  - Can use non-blocking IO
- Set FD to non-blocking using fcntl:
  - int x = fcntl(fd, F GETFL, 0);
  - $x \mid = O NONBLOCK;$
  - fcntl(fd, F SETFL, x);
- Now reads/writes/etc won't block
  - Just return immediately if can't perform IO immediately
  - Note: not magic

![](_page_49_Picture_11.jpeg)

• ONLY means that IO operation returns without waiting

![](_page_49_Picture_14.jpeg)

# **Non-Blocking IO: Continued**

```
int x = read (fd, buffer, size);
if (x < 0) {
   if (errno == EAGAIN) {
      //no data available
   else {
      //error
```

![](_page_50_Picture_2.jpeg)

![](_page_50_Picture_6.jpeg)

# **Non-Blocking IO: Continued**

```
while (size > 0) {
    int x = read (fd, buffer, size);
    if (x < 0) {
       if (errno == EAGAIN) {
          //no data available
       else {
          //error
    else {
      buffer += x;
      size -= x;
```

![](_page_51_Picture_2.jpeg)

#### What if we just wrap this up in a while loop?

![](_page_51_Picture_7.jpeg)

![](_page_51_Picture_8.jpeg)

![](_page_51_Picture_9.jpeg)

# **Non-Blocking IO: Continued**

```
while (size > 0) {
    int x = read (fd, buffer, size);
    if (x < 0) {
       if (errno == EAGAIN) {
          //no data available
       else {
          //error
    else {
      buffer += x;
      size -= x;
```

![](_page_52_Picture_2.jpeg)

- What if we just wrap this up in a while loop?
- Now we just made this blocking! We are just doing the blocking ourselves...

![](_page_52_Picture_8.jpeg)

![](_page_52_Picture_10.jpeg)

### • This approach is **worse** than blocking IO

• Why?

![](_page_53_Picture_3.jpeg)

![](_page_53_Picture_4.jpeg)

![](_page_53_Picture_6.jpeg)

# **Busy Wait**

- This approach is **worse** than blocking IO
  - Why?
- Busy waiting
  - Code is "actively" doing nothing
- Blocking IO:
  - At least OS will put thread to sleep while it waits

![](_page_54_Picture_8.jpeg)

• Keeping CPU busy, consuming power, contending with other threads

![](_page_54_Picture_14.jpeg)

### So What Do We Do?

### • Need to do **something else** while we wait

• Like what?

![](_page_55_Picture_3.jpeg)

![](_page_55_Picture_5.jpeg)

### So What Do We Do?

### • Need to do **something else** while we wait

- Like what?
- It depends....
  - On what?

![](_page_56_Picture_5.jpeg)

![](_page_56_Picture_7.jpeg)

# So What Do We Do?

- Need to do something else while we wait
  - Like what?
- It depends....
  - On what?
- On what our server does
- On what the demands on it are
- On the model of parallelism we are using

![](_page_57_Picture_9.jpeg)

• Who can name some models of parallelism? [AoP Ch 28 review]

![](_page_57_Picture_14.jpeg)

![](_page_58_Figure_1.jpeg)

- When would this be appropriate?
- What do our IO threads do for "something else"?

![](_page_58_Picture_4.jpeg)

![](_page_58_Picture_8.jpeg)

- When appropriate: Can keep IO thread(s) busy
  - Heavy IO to perform
  - Might have one thread do reads and writes
- What is "something else"?
  - Other IO requests
  - Do whichever one is ready to be done

![](_page_59_Picture_7.jpeg)

![](_page_59_Picture_12.jpeg)

- When appropriate: Can keep IO thread(s) busy
  - Heavy IO to perform
  - Might have one thread do reads and writes
- What is "something else"?
  - Other IO requests
  - Do whichever one is ready to be done
- Making hundreds of read/write calls to see which succeeds = inefficient
  - Use poll or select

![](_page_60_Picture_9.jpeg)

![](_page_60_Picture_11.jpeg)

![](_page_61_Figure_1.jpeg)

- What can you say about data movement in this model?
- What can you say about load balance?

![](_page_61_Picture_4.jpeg)

![](_page_61_Picture_8.jpeg)

![](_page_62_Picture_0.jpeg)

# Could have one thread work on many requests while(1) { Accept new requests Do any available reads/writes Do any available compute

- What can you say about data movement in this model?
- What can you say about load balance?

![](_page_62_Picture_4.jpeg)

# **Another Option**

![](_page_62_Picture_9.jpeg)

![](_page_63_Picture_0.jpeg)

### Slightly different inner loop: while(1) { Accept new requests For each request with anything to do Do any available IO for that request Do any compute for that request }

- What can you say about data movement in this model?
- What can you say about load balance?

![](_page_63_Picture_4.jpeg)

# **A Slight Variant**

![](_page_63_Picture_10.jpeg)