

Performance Tuning and Debugging

Don Porter

Why is my application slow?

- No silver bullet
- Part science, part art
 - Science: Measure performance, test hypotheses
 - Art: Finding practical balances of concerns

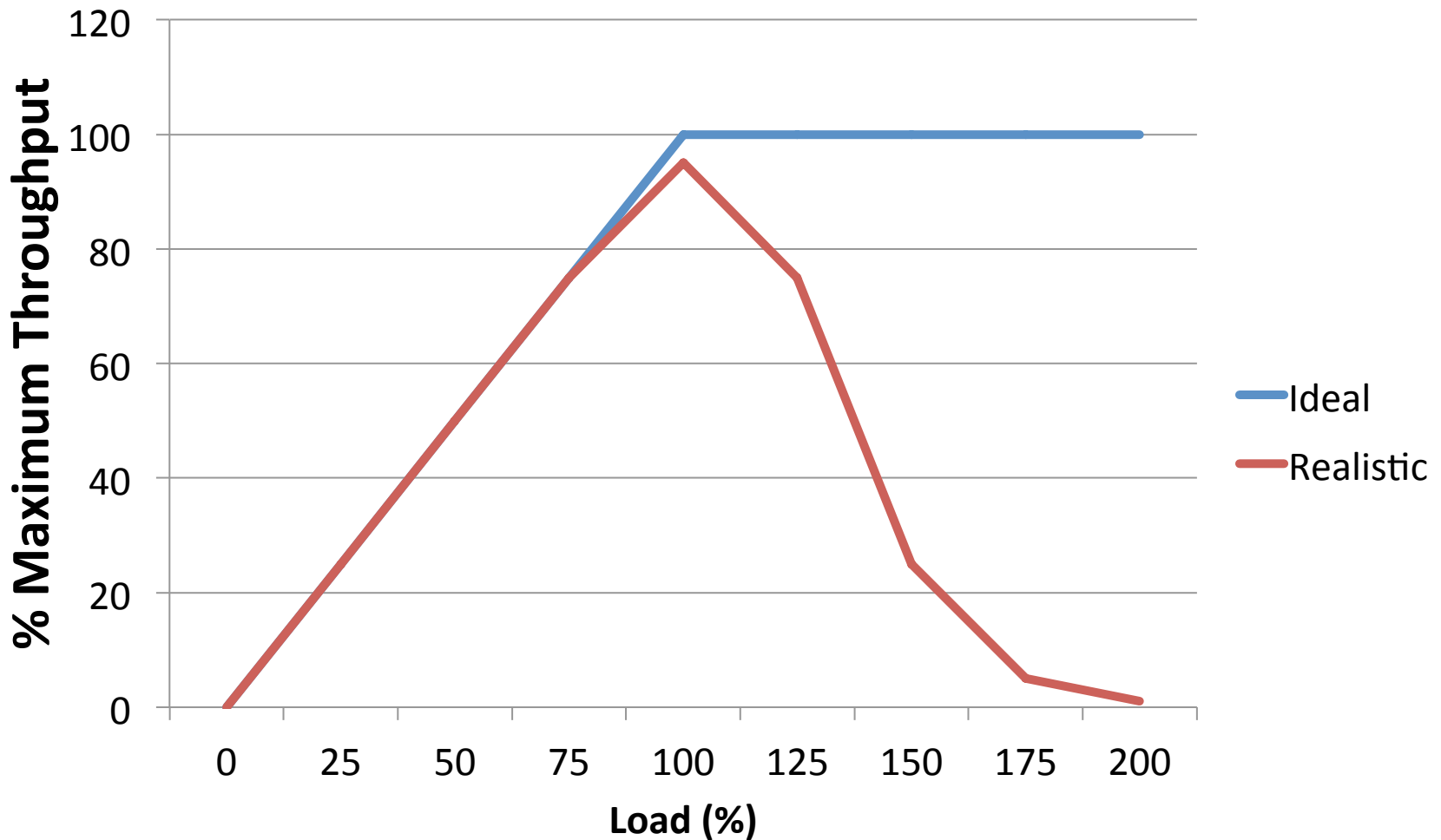
Most common culprits

- Insufficient resources
 - Configuration error
 - Hardware problems

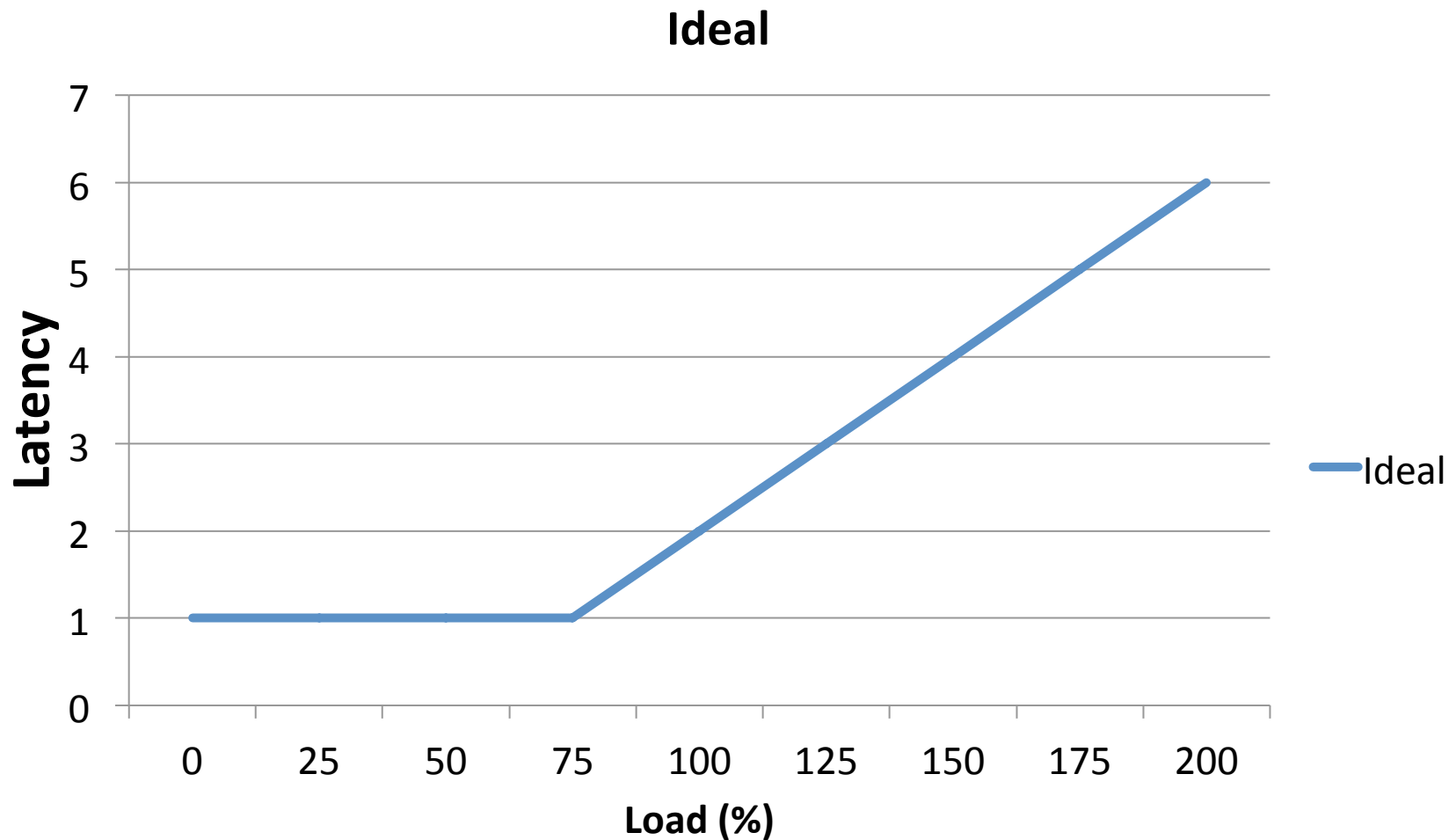
Digression: Throughput and Latency

- What are they?
- Throughput: Operations over time
 - Requests per second
 - Transactions per minute
 - Higher is better
- Latency: Time to complete one operation
 - My server can complete an HTTP GET in .01 seconds
 - Lower is better

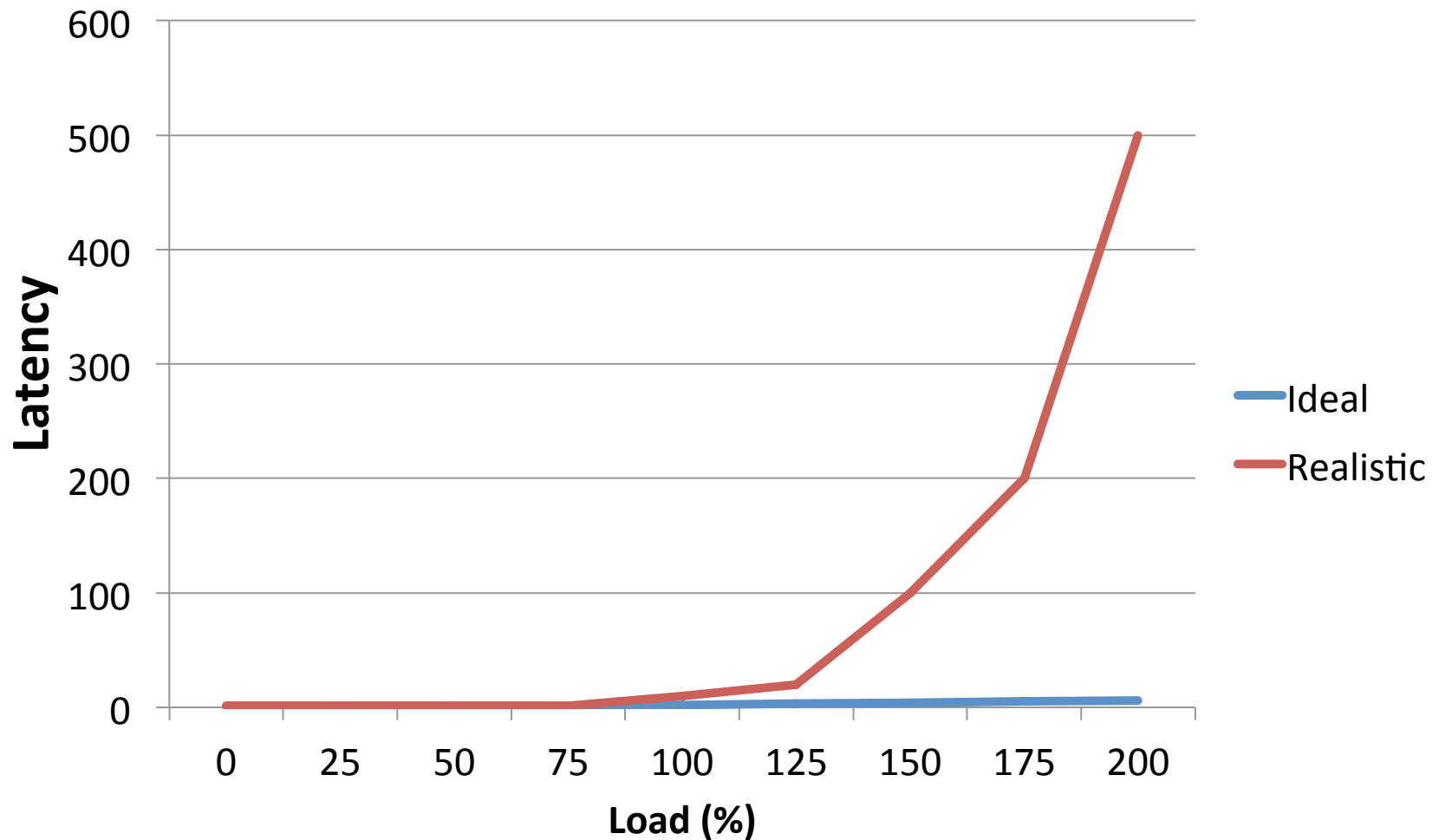
What happens when you are overloaded?



What happens when you are overloaded?



What happens when you are overloaded?



Note Change in Y Axis Scale---approaches infinity

Graceful Degradation

- Ideally, when a system is overloaded, by $n\%$, operation latency would increase by $n\%$ and throughput would stay constant
- In practice, systems rarely degrade gracefully when they are overloaded
- Thus, finding the “limiting factor” is essential

atop

ATOP - aria			2010/02/28 12:16:22			----			10s elapsed		
PRC	sys	5.20s	user	6.20s	#proc	157	#zombie	0	#exit	2	
CPU	sys	27%	user	61%	irq	25%	idle	214%	wait	73%	
cpu	sys	8%	user	59%	irq	25%	idle	8%	cpu003 w	0%	
cpu	sys	13%	user	1%	irq	0%	idle	85%	cpu001 w	0%	
cpu	sys	5%	user	1%	irq	0%	idle	22%	cpu002 w	72%	
cpu	sys	1%	user	0%	irq	0%	idle	99%	cpu000 w	0%	
CPL	avg1	2.08	avg5	1.91	avg15	1.31	csw	90799	intr	160344	
MEM	tot	7.6G	free	5.0G	cache	1.9G	buff	170.9M	slab	158.1M	
SWP	tot	2.0G	free	2.0G			vmcom	624.2M	vmlim	5.8G	
DSK	sda	busy	66%	read	0	write	1040	avio	6.36 ms		
DSK	sdb	busy	56%	read	0	write	848	avio	6.64 ms		
NET	transport		tcpi	188125	tcpo	99797	udpi	0	udpo	0	
NET	network		ipi	188125	ipo	99796	ipfrw	0	deliv	188125	
NET	eth0	21%	pcki	188120	pcko	99793	si	216 Mbps	so	5269 Kbps	

PID	SYSCPU	USRCPU	VGROW	RGROW	RDDSK	WRDSK	ST	EXC	S	CPUNR	CPU	CMD	1/2
17063	1.12s	4.98s	OK	OK	-	-	-E	0	E	-	61%	<bzip2>	
17059	2.25s	1.08s	OK	OK	OK	246.6M	--	-	R	3	33%	rsync	
17064	1.23s	0.13s	OK	OK	-	-	-E	0	E	-	14%	<tar>	
17011	0.27s	0.00s	OK	OK	OK	27012K	--	-	S	2	3%	pdflush	
16991	0.14s	0.00s	OK	OK	OK	10808K	--	-	S	2	1%	pdflush	
662	0.11s	0.00s	OK	OK	OK	4940K	--	-	S	2	1%	kjournald	
1957	0.04s	0.00s	OK	OK	OK	114.2M	--	-	S	1	0%	kjournald2	
17047	0.02s	0.01s	OK	OK	OK	OK	--	-	R	1	0%	atop	
2669	0.01s	0.00s	OK	OK	OK	OK	--	-	S	0	0%	httpd	
2045	0.01s	0.00s	OK	OK	OK	OK	--	-	S	0	0%	kondemand/0	
2687	0.00s	0.00s	OK	OK	OK	OK	--	-	S	3	0%	mythbackend	

atop

- Super-useful tool that shows usage of
 - CPU
 - Memory
 - Disk
 - Network
- On a color terminal, highlights over-used resources

CPU

- Very rarely the bottleneck
 - Actually degrades gracefully in most cases
- Nonetheless, overloaded CPUs will seem less responsive
- Note that when another resource is scarce, CPU time is used trying to compensate

Load Average

- The average number of processes waiting for the CPU
 - Less than 1, the CPU is idle
 - Higher than 1 is ok, just means CPU is fully utilized
 - Very high values (>8) can indicate a problem
- Read from the uptime command:

```
$ uptime
```

```
20:10:13 up 20 days, 11:08, 5 users, load average: 0.00, 0.03, 0.05
```

Memory

- Often the biggest troublemaker
- Why?
 - OSes over-commit memory to applications
 - In other words, if I have 1GB RAM, I can have 5 applications that all think they have 300 MB
 - How is this possible?
 - Swapping

Swapping

- If the OS is running low on memory, it can take RAM away from applications
 - Save the contents to disk
 - Reuse the RAM
- If the application tries to read or write to this memory, the application is interrupted, OS notified
 - OS has to then find free RAM, replace contents for app

The problem with swapping

- Disk reads and writes are slow (relative to CPU)
 - You very rarely wait for them before making progress
 - Except when swapping
- Mitigation: OS makes educated guesses about unlikely-to-be-used data to swap out
 - In the best case, things slow down a bit, and then return to normal
- In the worst case, data ping-pongs between disk and RAM
 - Called thrashing

Recommendation

- If you see substantial swap usage in atop, buy more RAM
 - It is cheap, and more RAM is cheaper now than when you bought the computer
- Note: OS often uses substantial amount of RAM to cache the file system contents, so don't be misled if total RAM usage is near 100%
 - Look at swap to detect insufficient RAM

In a crisis...

- Linux has an out-of-memory killer
- As advertised, it just kills programs until there is enough memory

Swappiness

- Linux tries to swap some data out before there is a crisis
- Linux has a parameter that sets how aggressively to swap data. This can get out of whack
 - `/proc/sys/vm/swappiness`
- I've personally had to dial this back on an Ubuntu release that set the default too high, in order for a nearly *idle* system to be usable

Network

- When the network is overloaded, packets are dropped
 - But the other end usually retries
- Two biggest culprit for network overload:
 - Attack (denial of service, brute-force password guessing, spam, etc)
 - Legitimate overload (slashdotted website, peak usage time)
- Need to figure out which

Network advice

- If the overload is not legitimate, good security practice can help to reduce wasteful traffic
 - Firewall, denyhosts, spam filter, etc.
 - For DoS, there are also quality-of-service tools on many network devices to limit the share of packets delivered from any one source
- If the overload is legitimate, you may need more servers and a load-balancer
 - Like round-robin DNS

Disks

- Very rarely the bottleneck, except:
 - (Implicitly when thrashing swap)
 - Actual disk-intensive workloads (e.g., database)
 - And when disk is nearing end-of-life
- Why rarely a problem?
 - Most disk requests are asynchronous
 - Most disk-intensive applications inherently rate-limited
- Why a problem at end-of-life?
 - Heavy remapping yields poor scheduling
 - For SSDs, internal bookkeeping can take longer as the device ages

Disks

- In general, if the disk is getting old, the best advice is replace it
 - You also don't want to lose data
- Some file systems perform worse as they age, but these are increasingly uncommon
 - Running a “defragmenter” can help

General advice

- Measure a performance baseline for your system
 - Application performance
 - Microbenchmarks (e.g., Imbench)
- If things seem slower, re-measure the component
 - Has my disk bandwidth degraded?
- This is the science of tuning

Other tools

- `/proc/cpuinfo`, `/proc/meminfo`, `/proc/diskstats` – useful system statistics
 - Lots of goodies in `/proc`
- `vmstat` – more details on memory usage
- `nice/renice` – adjust scheduling priority, giving more CPU time to important applications
- `swapinfo` – more details on swapping
- `netstat` – more details about network usage
- `hdparm/sdparm` – measure raw disk performance
- `iostat` – more details about disk I/O