

Going Beyond 100%*

The Little Hypervisor That Could

Valentin Bondzio

Sr. Staff TSE / Global Support

2021-07-20

“Light at the end of the tunnel” Edition

Biergarten

Vacations abroad

Eating out in peace

Licking each others' hands

Hugging your friends

DACH VMUG 2022 live?

Agenda

Intro

Motivation

CPU

Usage / Accounting

Memory

Consumed / Accounting


Takeaway*

*Hopefully

Motivation

Profitable + Sustainable = Efficient



A person wearing a yellow protective suit and a blue respirator mask with two circular filters is holding a smartphone. The person is giving a thumbs-up gesture. The background is a gradient of purple and pink.

“We are trying to stay around
60-70% average and not
exceed 85% peak utilization.”

Rule of Thumb

Navigation pane showing a tree view of vSphere objects. The root node is expanded to show a cluster, which is further expanded to show multiple hosts. The hosts are represented by small icons and labels, with some showing warning icons.

ACTIONS

Summary Monitor Configure Permissions **Hosts** VMs Datastores Networks Updates

Hosts Resource Pools

Name	State	Status	Cluster	Consumed CPU %	Consumed Memor...	HA State
...vmware.com	Connected	Warning	...	98%	...	? N/A
...vmware.com	Connected	Warning	...	98%	...	? N/A
...vmware.com	Connected	Warning	...	98%	...	? N/A
...vmware.com	Connected	Warning	...	98%	...	? N/A
...vmware.com	Connected	Warning	...	98%	...	? N/A
...vmware.com	Connected	Warning	...	98%	...	? N/A
...vmware.com	Connected	Warning	...	98%	...	? N/A
...vmware.com	Connected	Warning	...	98%	...	? N/A
...vmware.com	Connected	Warning	...	98%	...	? N/A
...vmware.com	Connected	Warning	...	98%	...	? N/A

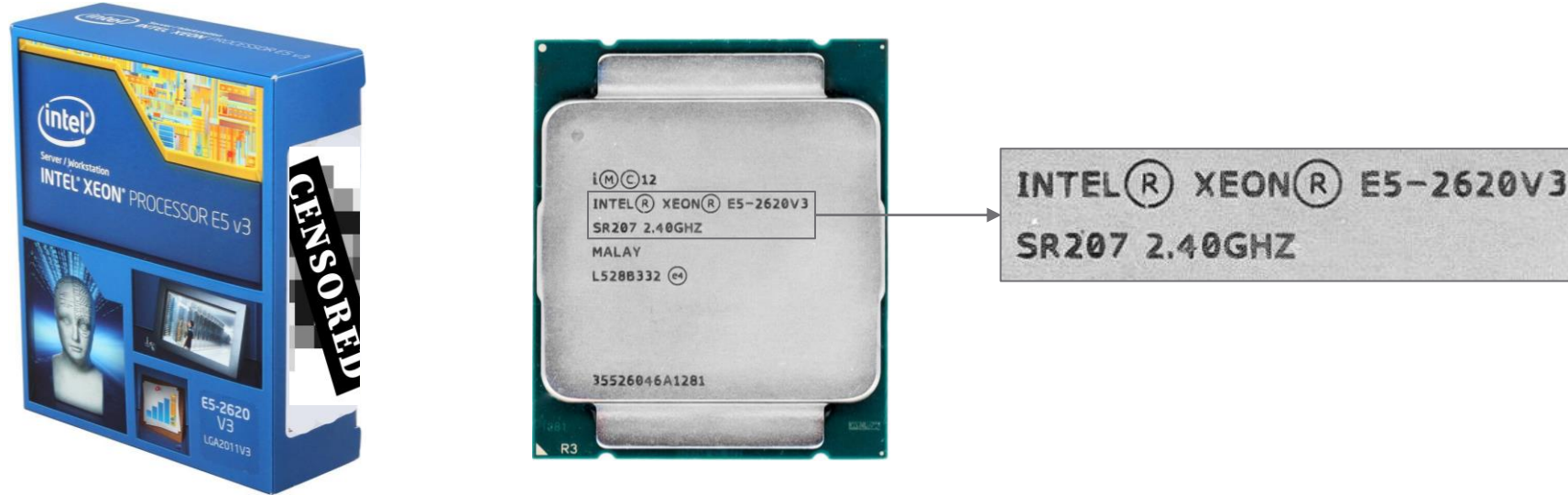


CPU Usage (host level) is accounted:
at **nominal frequency** and **without SMT**

100% Host CPU Usage =
nominal frequency * number of cores

Nominal (Base) Frequency

“What’s on the box”



```
[root@cs-tse-d93:~] vsish -e get /hardware/cpu/cpuModelName  
Intel(R) Xeon(R) CPU E5-2620 v3 @ 2.40GHz
```

```
[root@cs-tse-d93:~] gzip -dc /var/log/boot.gz | grep -i timer.*measured  
TSC: 4945939461 cpu0:1)Timer: 737: cpu 0: early measured tsc speed is 2399997831 Hz  
0:00:00:02.412 cpu0:1)Timer: 1196: cpu 0: measured cpu speed: 2399998045 Hz, TSC speed: 2399997831 Hz, bus speed: 99999904 Hz
```

Nominal (Base) Frequency and Beyond

Turbo Boost



Table 2. Intel® Xeon® Processor E5-1600, E5-2400, E5-2600 and E5-4600 v3 Product Families Turbo Bins (Sheet 2 of 3)

S-Spec No	Stepping	Model Number	TDP (W)	# Cores	Intel® Turbo Boost Technology Maximum Core Frequency (GHz)											Notes
					Core 1-2	Core 3	Core 4	Core 5	Core 6	Core 7	Core 8	Core 9	Core 10	Core 11+		
SR1XV	M1	E5-2658 v3	105	12	2.9	2.7	2.6	2.5	2.5	2.5	2.5	2.5	2.5	2.5	1,2,3,7	
SR1XW	M1	E5-2648L v3	75	12	2.5	2.3	2.2	2.1	2.1	2.1	2.1	2.1	2.1	2.1	1,2,3,7	
SR1XZ	M1	E5-2628L v3	75	10	2.5	2.3	2.2	2.2	2.2	2.2	2.2	2.2	2.2	NA	1,2,3,7	
SR20H	R2	E5-1680 v3	140	8	3.8	3.6	3.5	3.5	3.5	3.5	3.5	NA	NA	NA	1,2,3,6,7	
SR20J	R2	E5-1650 v3	140	6	3.8	3.6	3.6	3.6	3.6	NA	NA	NA	NA	NA	1,2,3,6,7	
SR207	R2	E5-2620 v3	85	6	3.2	2.9	2.8	2.7	2.6	NA	NA	NA	NA	NA	1,2,3,7	
SR1YC	M1	E5-2609 v3	85	6	1.9	1.9	1.9	1.9	1.9	NA	NA	NA	NA	NA	1,2,3,4,5	
SR20N	R2	E5-1660 v3	140	8	3.5	3.3	3.3	3.3	3.3	3.3	3.3	NA	NA	NA	1,2,3,6,7	
SR20B	R2	E5-2608L v3	50	6	2	2	2	2	2	NA	NA	NA	NA	NA	1,2,3,7	
SR20P	R2	E5-1620 v3	140	4	3.6	3.6	3.6	NA	NA	NA	NA	NA	NA	NA	1,2,3,6,7	
SR21P	R2	E5-2608L v3	52	6	2	2	2	2	2	NA	NA	NA	NA	NA	1,2,3,7	
SR22P	C1	E5-4660 v3	120	14	2.9	2.7	2.6	2.5	2.4	2.4	2.4	2.4	2.4	2.4	1,2,3,7,8	
SR22J	C1	E5-4650 v3	105	12	2.8	2.6	2.5	2.4	2.4	2.4	2.4	2.4	2.4	2.4	1,2,3,7,8	
SR22L	C1	E5-4640 v3	105	12	2.6	2.4	2.3	2.2	2.1	2.1	2.1	2.1	2.1	2.1	1,2,3,7,8	

```
[root@cs-tse-d93:~] sched-stats -t ncpus
24 PCPUs
12 cores
2 LLCs
2 packages
2 NUMA nodes
```

```
[root@cs-tse-d93:~] esxtop
8:49:09pm up 52 days 11:49, 970 worlds, 0 VMs, 0 vCPUs; (...)
Power Usage: 140W, Power Cap: N/A
PSTATE MHZ: 2401 2400 2300 (...) 1300 1200
```

```
CPU %USED %UTIL %C0 %C1 %C2 %P0 %P1 %P2 (...) %P12 %P13 %A/MPERF
 0 0.0 0.1 0 100 0 100 0 0 (...) 0 0 108.3
 1 0.0 0.0 0 100 0 100 0 0 (...) 0 0 108.4
 (...)
11 0.1 0.1 0 100 0 100 0 0 (...) 0 0 108.3
12 0.1 0.1 0 100 0 100 0 0 (...) 0 0 108.3
 (...)
22 0.0 0.0 0 100 0 100 0 0 (...) 0 0 108.2
23 0.0 0.0 0 100 0 100 0 0 (...) 0 0 108.3
```

```
[root@cs-tse-d93:~] awk "BEGIN{printf \"%.1f\\n\", (26 / 24) * 100}"
108.3
```

Nominal (Base) Frequency and Beyond

Turbo Boost



Table 2. Intel® Xeon® Processor E5-1600, E5-2400, E5-2600 and E5-4600 v3 Product Families Turbo Bins (Sheet 2 of 3)

S-Spec No	Stepping	Model Number	TDP (W)	# Cores	Intel® Turbo Boost Technology Maximum Core Frequency (GHz)										Notes
					Core 1-2	Core 3	Core 4	Core 5	Core 6	Core 7	Core 8	Core 9	Core 10	Core 11+	
SR1XV	M1	E5-2658 v3	105	12	2.9	2.7	2.6	2.5	2.5	2.5	2.5	2.5	2.5	2.5	1,2,3,7
SR1XW	M1	E5-2648L v3	75	12	2.5	2.3	2.2	2.1	2.1	2.1	2.1	2.1	2.1	2.1	1,2,3,7
SR1XZ	M1	E5-2628L v3	75	10	2.5	2.3	2.2	2.2	2.2	2.2	2.2	2.2	2.2	NA	1,2,3,7
SR20H	R2	E5-1680 v3	140	8	3.8	3.6	3.5	3.5	3.5	3.5	3.5	NA	NA	NA	1,2,3,6,7
SR20J	R2	E5-1650 v3	140	6	3.8	3.6	3.6	3.6	3.6	NA	NA	NA	NA	NA	1,2,3,6,7
SR207	R2	E5-2620 v3	85	6	3.2	2.9	2.8	2.7	2.6	NA	NA	NA	NA	NA	1,2,3,7
SR1YC	M1	E5-2609 v3	85	6	1.9	1.9	1.9	1.9	1.9	NA	NA	NA	NA	NA	1,2,3,4,5
SR20N	R2	E5-1660 v3	140	8	3.5	3.3	3.3	3.3	3.3	3.3	3.3	NA	NA	NA	1,2,3,6,7
SR20B	R2	E5-2608L v3	50	6	2	2	2	2	2	NA	NA	NA	NA	NA	1,2,3,7
SR20P	R2	E5-1620 v3	140	4	3.6	3.6	3.6	NA	NA	NA	NA	NA	NA	NA	1,2,3,6,7
SR21P	R2	E5-2608L v3	52	6	2	2	2	2	2	NA	NA	NA	NA	NA	1,2,3,7
SR22P	C1	E5-4660 v3	120	14	2.9	2.7	2.6	2.5	2.4	2.4	2.4	2.4	2.4	2.4	1,2,3,7,8
SR22J	C1	E5-4650 v3	105	12	2.8	2.6	2.5	2.4	2.4	2.4	2.4	2.4	2.4	2.4	1,2,3,7,8
SR22L	C1	E5-4640 v3	105	12	2.6	2.4	2.3	2.2	2.1	2.1	2.1	2.1	2.1	2.1	1,2,3,7,8

```
[root@cs-tse-d93:~] sched-stats -t ncpus
24 PCPUs
12 cores
2 LLCs
2 packages
2 NUMA nodes
```

```
[root@cs-tse-d93:~] esxtop
8:49:09pm up 52 days 11:49, 970 worlds, 0 VMs, 0 vCPUs; (...)
Power Usage: 140W, Power Cap: N/A
PSTATE MHZ: 2401 2400 2300 (...) 1300 1200
```

```
CPU %USED %UTIL %C0 %C1 %C2 %P0 %P1 %P2 (...) %P12 %P13 %A/MPERF
 0 0.0 0.1 0 100 0 100 0 0 (...) 0 0 108.3
 1 0.0 0.0 0 100 0 100 0 0 (...) 0 0 108.4
 (...)
11 0.1 0.1 0 100 0 100 0 0 (...) 0 0 108.3
12 0.1 0.1 0 100 0 100 0 0 (...) 0 0 108.3
 (...)
22 0.0 0.0 0 100 0 100 0 0 (...) 0 0 108.2
23 0.0 0.0 0 100 0 100 0 0 (...) 0 0 108.3
```

```
[root@cs-tse-d93:~] awk "BEGIN{printf \"%.1f\\n\", (26 / 24) * 100}"
108.3
```


Identifying MACF

Not a common industry acronym!

Configure whatever “Maximum / High Performance” Power Management policy available

[root@cs-tse-d93:~] esxtop
(...)

p, f, f

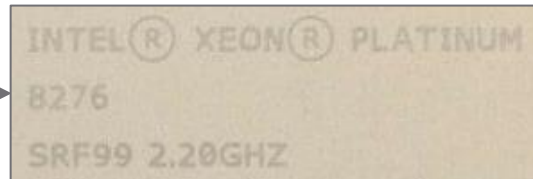
4:38:51pm up 1 min, 1276 worlds, 0 VMs, 0 vCPUs; CPU load average: 0.02, 0.00, 0.00
Power Usage: 142W, Power Cap: N/A
PSTATE MHZ:

CPU	%USED	%UTIL	%C0	%C1	%A/MPERF
0	0.0	0.1	0	100	108.3
1	0.1	0.1	0	100	108.4
2	0.1	0.1	0	100	108.3
3	0.0	0.1	0	100	108.4
(...)					
20	0.0	0.1	0	100	108.4
21	0.0	0.0	0	100	108.3
22	0.0	0.0	0	100	108.3
23	0.2	0.2	0	100	108.3

MACF is 8.3 % on top of NF for this CPU / Host

Cycle “Capacity”

Intel Xeon Platinum 8276



2.20 GHz NF
 3.00 GHz MACF
 4.00 GHz M2CF

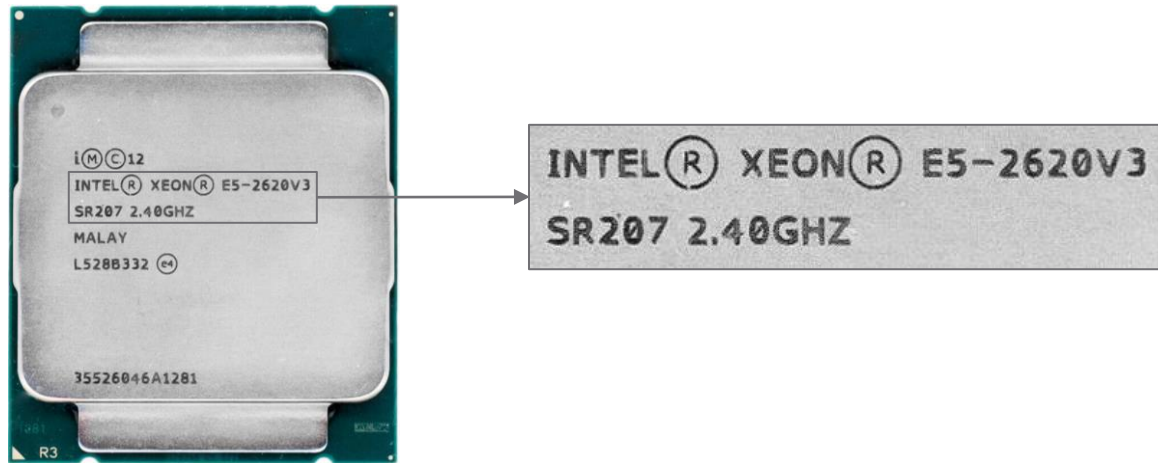
36% additional capacity @ MACF
 82% additional speed @ M2CF*

* needs deep C-States / non “High Performance” policy

Mode	Base	Turbo Frequency/Active Cores																												
		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	
Normal	2,200MHz	4,000MHz	4,000MHz	3,800MHz	3,800MHz	3,700MHz	3,700MHz	3,700MHz	3,700MHz	3,700MHz	3,700MHz	3,700MHz	3,700MHz	3,700MHz	3,700MHz	3,700MHz	3,700MHz	3,400MHz	3,400MHz	3,400MHz	3,400MHz	3,100MHz	3,100MHz	3,100MHz	3,100MHz	3,000MHz	3,000MHz	3,000MHz	3,000MHz	
AVX2	1,700MHz	3,800MHz	3,800MHz	3,600MHz	3,600MHz	3,500MHz	3,500MHz	3,500MHz	3,500MHz	3,500MHz	3,500MHz	3,500MHz	3,500MHz	3,200MHz	3,200MHz	3,200MHz	3,200MHz	2,900MHz	2,900MHz	2,900MHz	2,900MHz	2,700MHz	2,700MHz	2,700MHz	2,700MHz	2,600MHz	2,600MHz	2,600MHz	2,600MHz	
AVX512	1,300MHz	3,700MHz	3,700MHz	3,500MHz	3,500MHz	3,300MHz	3,300MHz	3,300MHz	3,300MHz	3,300MHz	2,900MHz	2,900MHz	2,900MHz	2,900MHz	2,600MHz	2,600MHz	2,600MHz	2,600MHz	2,300MHz	2,300MHz	2,300MHz	2,300MHz	2,200MHz	2,200MHz	2,200MHz	2,200MHz	2,100MHz	2,100MHz	2,100MHz	2,100MHz

Cycle “Capacity”

Intel Xeon E5-2620 v3



2.40 GHz NF
 2.60 GHz MACF
 3.20 GHz M2CF

8% additional capacity @ MACF
 33% additional speed @ M2CF*

* needs deep C-States / non “High Performance” policy



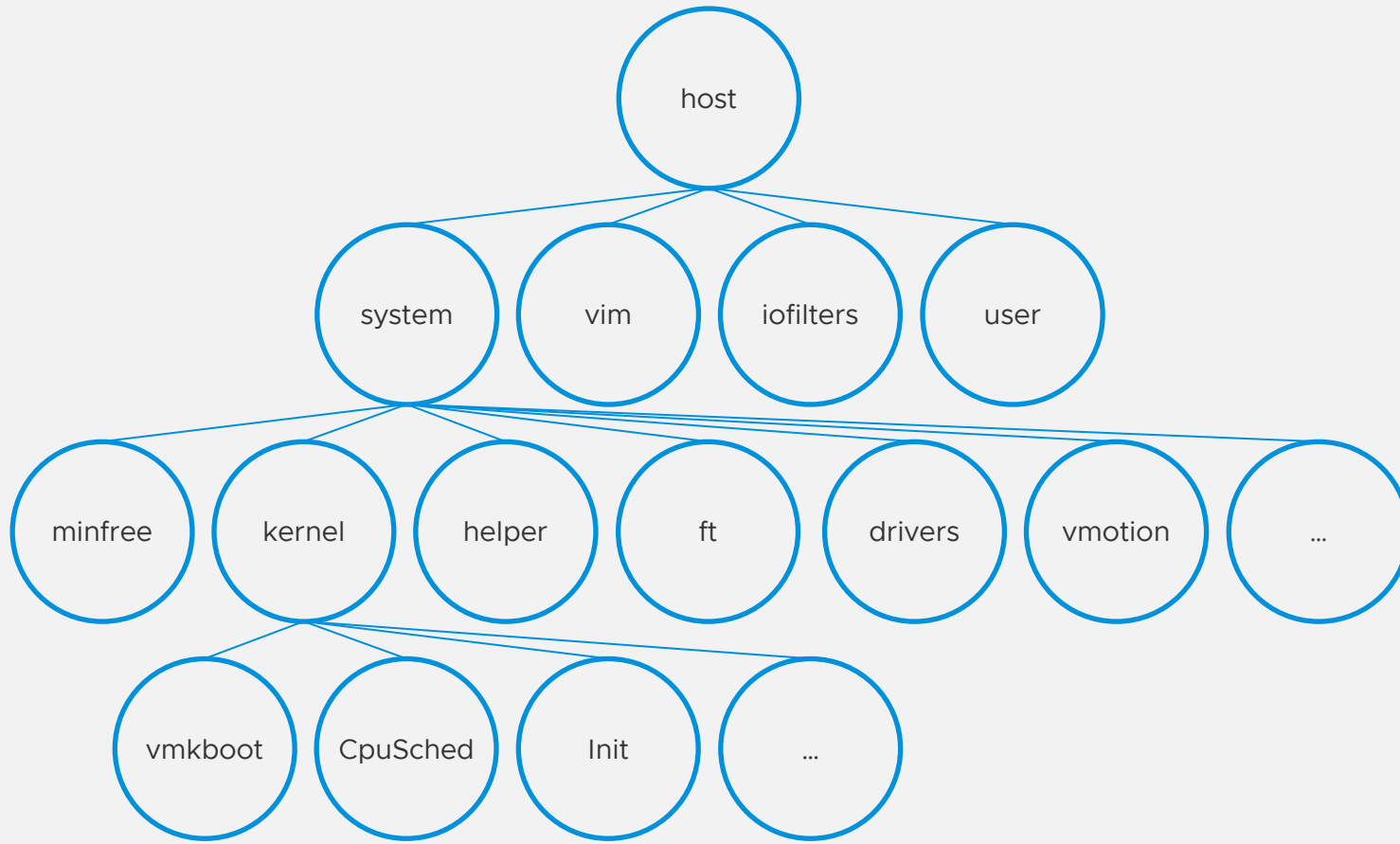
Table 2. Intel® Xeon® Processor E5-1600, E5-2400, E5-2600 and E5-4600 v3 Product Families Turbo Bins (Sheet 2 of 3)

S-Spec No	Stepping	Model Number	TDP (W)	# Cores	Intel® Turbo Boost Technology Maximum Core Frequency (GHz)										Notes
					Core 1-2	Core 3	Core 4	Core 5	Core 6	Core 7	Core 8	Core 9	Core 10	Core 11+	
SR207	R2	E5-2620 v3	85	6	3.2	2.9	2.8	2.7	2.6	NA	NA	NA	NA	NA	1,2,3,7

Reservable == Capacity?

Hierarchical Resource Groups

From an ESXi perspective



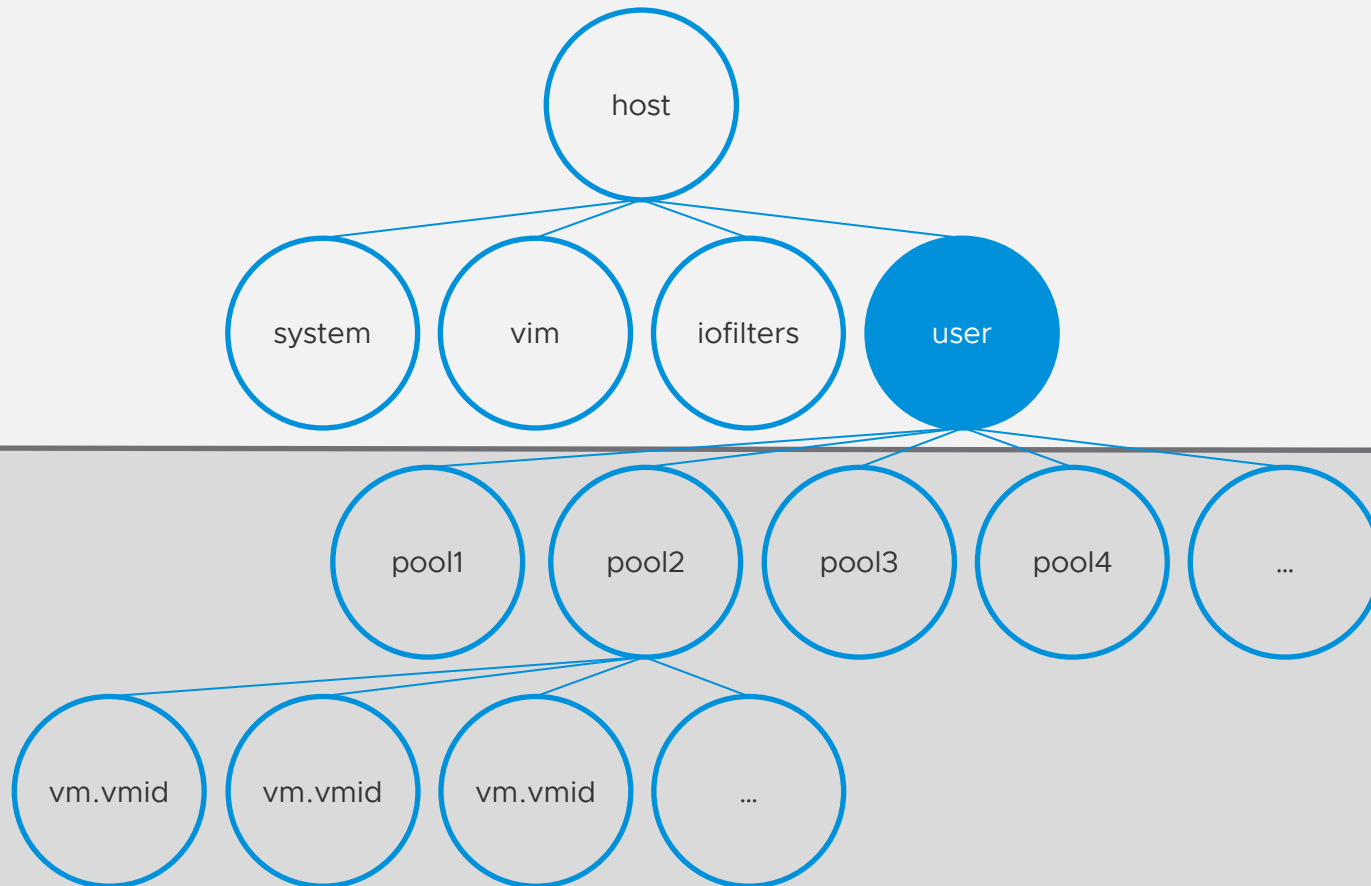
The host owns all resources

Those are distributed by hierarchical resource groups

Consumers can demand (request) resources

Hierarchical Resource Groups

From an ESXi perspective



vCenter shows the sum of all **user** resources as:

Total Reservation Capacity

Global Resource Pools are then distributed back to hosts into Local RPs

- Based on VMs demand

Cycle “Capacity”

Intel Xeon E5-2620 v3

```
[root@cs-tse-d93:~] vsish -e get /sched/groups/0/stats/capacity
group-capacity {
  cpu-reserved:4296 MHz
  cpu-unreserved:24504 MHz
  mem-reserved:7005904 KB
  mem-unreserved:59674248 KB
}
[root@cs-tse-d93:~] awk "BEGIN{printf \"%.0i\\n\", (4296 + 24504) / 12 }"
2400
```

```
[root@cs-tse-d93:~] vsish -e get /sched/groups/4/stats/capacity
group-capacity {
  cpu-reserved:0 MHz
  cpu-unreserved:24504 MHz
  mem-reserved:33340 KB
  mem-unreserved:59674248 KB
}
```

```
[root@cs-tse-d93:~] sched-stats -t groups -z vmgid:name:pgid:pname:size:vsmps:amin:amax:units:resvMHz:availMHz
| awk 'NR == 1 || $2 ~ /^(vm\.|pool)[0-9]+/ || /^[0-4] / {print $0}'
```

vmgid	name	pgid	pname	size	vsmps	amin	amax	units	resvMHz	availMHz
0	host	0	none	4	969	1200	1200	pct	4296	24504
1	system	0	host	690	695	10	-1	pct	288	24504
2	vim	0	host	4	271	4008	-1	mhz	3648	24864
3	iofilters	0	host	3	3	0	-1	pct	0	24504
4	user	0	host	0	0	0	-1	pct	0	24504

-z in 7.0.2

System CPU Usage

2nd level resource pools

Chart Options | cs-tse-d93.csl.vmware.com

Chart options: --Select option-- Save Options As... Delete Options

Chart Metrics

- CPU
- Cluster services
- Datstore
- Disk
- Memory
- Network
- Power
- Storage adapter
- Storage path
- System**
- vSphere Replication

Select counters for this chart:

Counters	Rollups	Units	Internal Name	Stat Type	Description	
<input type="checkbox"/>	Resource CPU allocation minimum (in MHz)	Latest	MHz	resourceCpuAllocMin	Absolute	CPU allocation r...
<input type="checkbox"/>	Resource CPU allocation shares	Latest	num	resourceCpuAllocShares	Absolute	CPU allocation s...
<input type="checkbox"/>	Resource CPU maximum limited (1 min)	Latest	%	resourceCpuMaxLimited1	Absolute	CPU maximum li...
<input type="checkbox"/>	Resource CPU maximum limited (5 min)	Latest	%	resourceCpuMaxLimited5	Absolute	CPU maximum li...
<input type="checkbox"/>	Resource CPU running (1 min. average)	Latest	%	resourceCpuRun1	Absolute	CPU running eve...
<input type="checkbox"/>	Resource CPU running (5 min average)	Latest	%	resourceCpuRun5	Absolute	CPU running eve...
<input checked="" type="checkbox"/>	Resource CPU usage (Average)	Average	MHz	resourceCpuUsage	Rate	Amount of CPU ...
<input type="checkbox"/>	Resource memory allocation maximum (in KB)	Latest	KB	resourceMemAllocMax	Absolute	Memory allocatio...
<input type="checkbox"/>	Resource memory allocation minimum (in KB)	Latest	KB	resourceMemAllocMin	Absolute	Memory allocatio...
<input type="checkbox"/>	Resource memory allocation shares	Latest	num	resourceMemAllocShares	Absolute	Memory allocatio...

Timespan: **Real-time**

Last: 1 Hour(s)

From: 24 Aug 2020 12:44:40

To: 25 Aug 2020 12:44:40
(time is in ISO 8601 format)

Chart Type: **Line Graph**

Select object for this chart:

- Target Objects
- host
- host/lofilters
- host/lofilters/lofiltervpd
- host/lofilters/spm
- host/lofilters/vmwarevmcrypt
- host/system
- host/system/drivers
- host/system/ft
- host/system/helper
- host/system/lfhelper

cs-tse-d93.csl.vmware.com | ACTIONS

Summary Monitor Configure Permissions VMs Datastores Networks Updates

Advanced Performance



System, 08/25/2020, 11:48:40 AM - 08/25/2020, 12:48:20 PM Real-time Chart Options View: Custom

Performance Chart Legend

Key	Object	Measurement	Rollup	Units	Latest	Maximum	Minimum	Average
■	host/lofilters	Resource CPU usage (Average)	Average	MHz	1	2	0	0.606
■	host/user	Resource CPU usage (Average)	Average	MHz	12	18	7	12.017
■	host/vim	Resource CPU usage (Average)	Average	MHz	23	414	8	21.072
■	host/system	Resource CPU usage (Average)	Average	MHz	16	300	6	50.017

Test Setup

Scripts and Helpers

 **esxi_one-host-sized-vm_settings-echo.sh**  769 Bytes

```
1 schedNcpus=$(sched-stats -t ncpus)
2 numPcpus=$(echo "${schedNcpus}" | sed -n 's/\([0-9]\+\) PCPUs$/\1/p')
3 numCores=$(echo "${schedNcpus}" | sed -n 's/\([0-9]\+\) cores$/\1/p')
4 numNumaNodes=$(echo "${schedNcpus}" | sed -n 's/\([0-9]\+\) NUMA nodes$/\1/p')
5 # don't care for package size / LLC etc.
6 schedDomainSize=$(( ${numPcpus} / ${numNumaNodes} ))
7
8 echo "cpuid.coresPerSocket = ${schedDomainSize}"
9 echo "numa.vcpu.maxPerVirtualNode = ${schedDomainSize}"
10 # theoretically only necessary with > 2 NUMA nodes but lets be safe
11 # [ "${numPcpus} / ${numCores}" -gt "1" -a "${numNumaNodes}" -gt "2" ]
12 if [ "${numPcpus} / ${numCores}" -gt "1" ]; then
13     echo "numa.vcpu.preferHT = true"
14 fi
15 for i in $(seq 0 1 ${numPcpus})
16 do echo -e "sched.vcpu${i}.affinity = ${i}"
17 done
```

```
# ./esxi_one-host-sized-vm_settings-echo.sh
```

```
cpuid.coresPerSocket = 12
numa.vcpu.maxPerVirtualNode = 12
numa.vcpu.preferHT = true
sched.vcpu0.affinity = 0
sched.vcpu1.affinity = 1
sched.vcpu2.affinity = 2
sched.vcpu3.affinity = 3
sched.vcpu4.affinity = 4
sched.vcpu5.affinity = 5
sched.vcpu6.affinity = 6
sched.vcpu7.affinity = 7
(...)
sched.vcpu16.affinity = 16
sched.vcpu17.affinity = 17
sched.vcpu18.affinity = 18
sched.vcpu19.affinity = 19
sched.vcpu20.affinity = 20
sched.vcpu21.affinity = 21
sched.vcpu22.affinity = 22
sched.vcpu23.affinity = 23
sched.vcpu24.affinity = 24
```

Test Setup

Scripts and Helpers

esxi_pre-bench_cosched-control.sh 496 Bytes

```
1  if [ -z "$1" ]; then
2      echo "call with either \"disable\", \"max\" or \"reset\" (to default)"; exit 0
3  fi
4
5  case $1 in
6
7      "disable")
8          start=0
9          stop=0
10         echo "disabled cosched"
11         ;;
12
13     "max")
14         start=99000
15         stop=100000
16         echo "maxed cosched"
17         ;;
18
19     "reset")
20         start=2000
21         stop=3000
22         echo "reset cosched to system defaults"
23         ;;
24
25  esac
26
27  vsish -e set /config/Cpu/intOpts/CoschedCostartThreshold ${start} &>/dev/null
28  vsish -e set /config/Cpu/intOpts/CoschedCostopThreshold ${stop} &>/dev/null
```

```
# ./esxi_pre-bench_cosched-control.sh disable
disabled cosched
```

Test Setup

Scripts and Helpers

```
esxi_pre-bench_power-stabilizer.sh 992 Bytes Edit Web IDE
1  if [ -z "$1" ]; then
2      echo "call with either \"set\" or \"reset\" (to default)"; exit 0
3  fi
4
5  hwSupport=$(vsish -e get /power/hardwareSupport | sed -n 's/ CPU power management:\(ACPI P-states, ACPI C-states\)/\1/p')
6
7  if [ -z "$hwSupport" ]; then
8      echo "Check that BIOS is properly configured"; exit 0
9  fi
10
11  cpuSupport=$(vsish -e get /hardware/cpu/cpuList/0 | sed -n 's/ Name:\(GenuineIntel\)/\1/p')
12
13  if [ -z "$cpuSupport" ]; then
14      echo "Only known to work on Intel"; exit 0
15  fi
16
17  case $1 in
18
19      "set")
20          policy=4
21          max=99
22          min=99
23          dc=0
24          echo "disabled Turbo Boost and deep C-States, vote for P1/NF"
25          ;;
26
27      "reset")
28          policy=2
29          max=100
30          min=0
31          dc=1
32          echo "reset custom power options to system default and policy to balanced"
33          ;;
34
35  esac
36
37  vsish -e set /power/currentPolicy ${policy}
38  vsish -e set /config/Power/intOpts/MaxFreqPct ${max} &>/dev/null
39  vsish -e set /config/Power/intOpts/MinFreqPct ${min} &>/dev/null
40  vsish -e set /config/Power/intOpts/UseCStates ${dc} &>/dev/null
```

```
# ./esxi_pre-bench_power-stabilizer.sh set
disabled Turbo Boost and deep C-States, vote for P1/NF
```

Test Setup

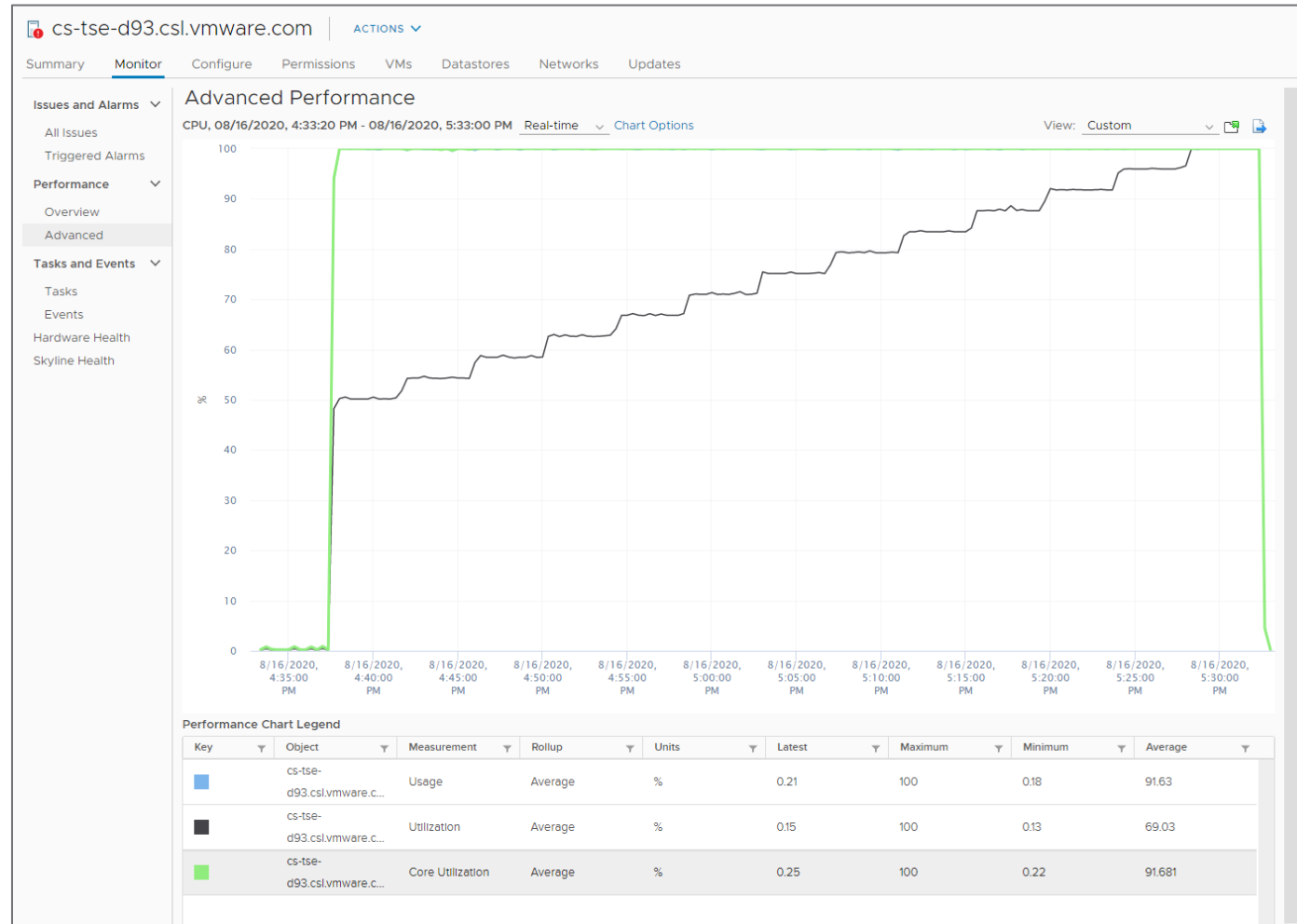
Scripts and Helpers

```
linux_stress-ng_usage-util_loadramp.sh 1007 Bytes Edit Web IDE Lock Replace D
1 if [ -z "$1" ]; then
2     echo "call with runtime in minutes (default: 5)"
3     # I usually select a duration that fits into realtime charts, so max. 60
4     testDurationMinutes=5
5 else
6     testDurationMinutes=$1
7 fi
8
9 # maybe include --cpu-load and change --cpu-method from loop to something else
10 cpuLoadMethod="loop"
11
12 cpus=$(nproc --all)
13 testDurationSeconds=$((testDurationMinutes*60))
14 testIntervalSeconds=$((testDurationSeconds / ((cpus / 2) + 1)))
15 initialTaskset=$(seq 0 2 $((cpus - 1)) | paste -s -d ,)
16
17 # run stress-ng in background and quiet
18 # baseline -> half of the vCPUs, all of the cores
19 stress-ng --taskset $initialTaskset --cpu $((cpus / 2)) --cpu-method $cpuLoadMethod -t $testDurationSeconds -q &
20
21 loop=0
22 for iterateTaskset in $(seq 1 2 $((cpus - 1)))
23 do
24     loop=$((loop + 1))
25     sleep $testIntervalSeconds
26     # add another worker to each hypervin
27     stress-ng --taskset $iterateTaskset --cpu 1 --cpu-method $cpuLoadMethod -t $((testDurationSeconds - (loop * testIntervalSeconds))) -q &
28 done
```

```
# ./ linux_stress-ng_usage-util_loadramp.sh
stress-ng: info: [2293] dispatching hogs: 12 cpu
stress-ng: info: [2311] dispatching hogs: 1 cpu
stress-ng: info: [2317] dispatching hogs: 1 cpu
stress-ng: info: [2323] dispatching hogs: 1 cpu
stress-ng: info: [2329] dispatching hogs: 1 cpu
stress-ng: info: [2335] dispatching hogs: 1 cpu
stress-ng: info: [2341] dispatching hogs: 1 cpu
stress-ng: info: [2347] dispatching hogs: 1 cpu
stress-ng: info: [2366] dispatching hogs: 1 cpu
stress-ng: info: [2372] dispatching hogs: 1 cpu
stress-ng: info: [2387] dispatching hogs: 1 cpu
stress-ng: info: [2399] dispatching hogs: 1 cpu
stress-ng: info: [2448] dispatching hogs: 1 cpu
```


Host view

VM starts with 12 vCPUs / all cores utilized



CONFIDENTIAL

```
hostd_getting_stats_from_vmksnel.csomething - Notepad
File Edit Format View Help
int MAXPCT = 100;
( ... )
AddFancySample(many secret arguments, time, cpuUsage, MAXPCT);
Ln 4, Col 1 100% Windows (CRLF) UTF-8
```

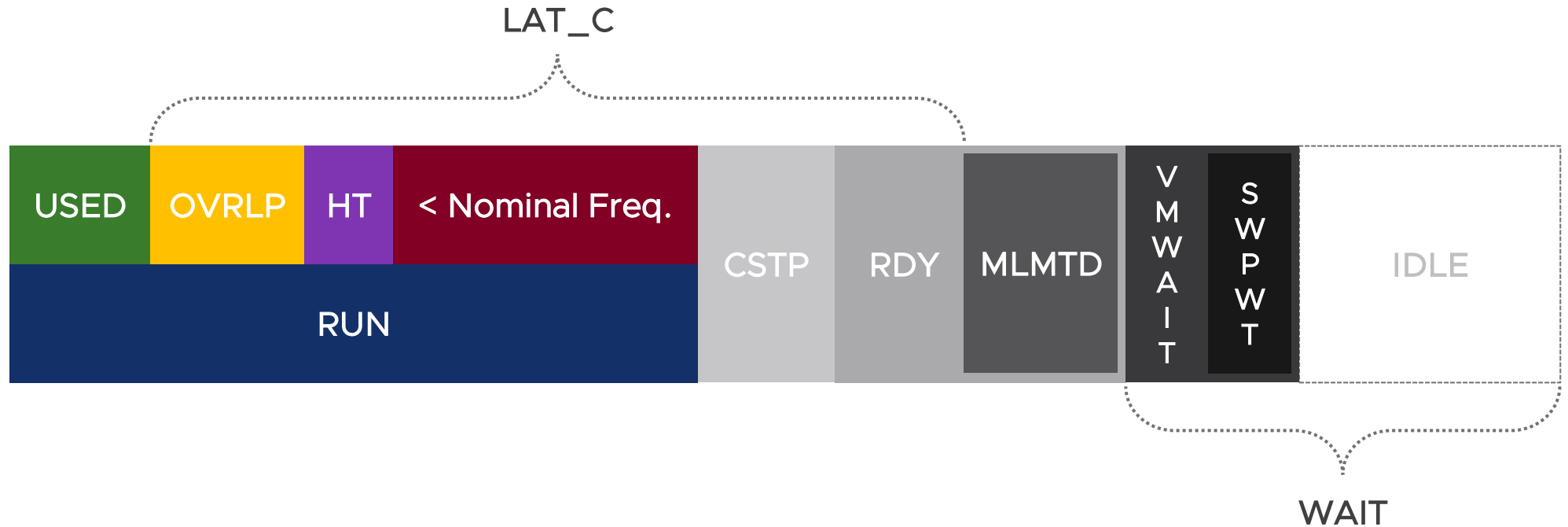
vCPU accounting

States and Metrics – esxtop nomenclature



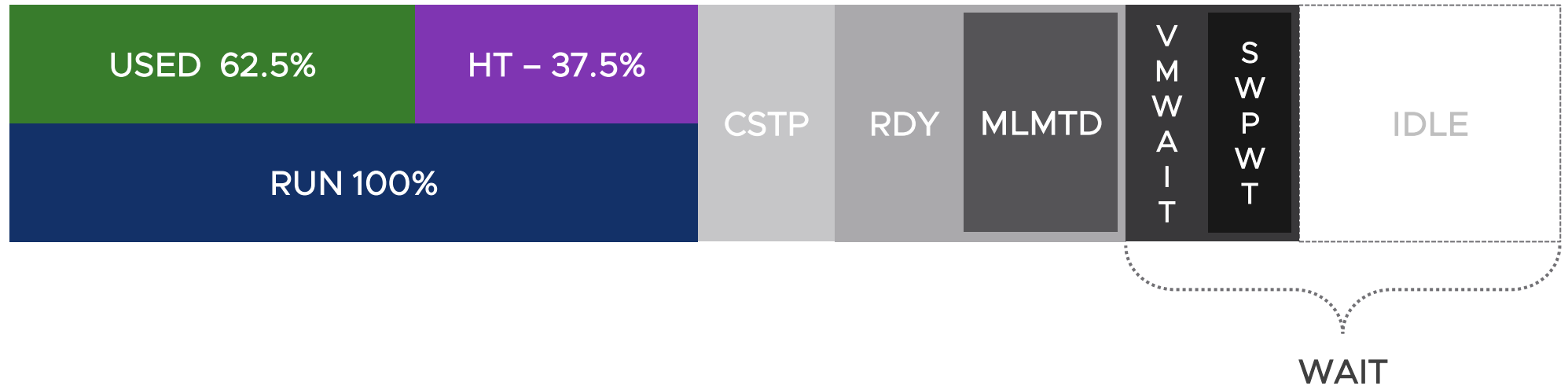
vCPU accounting

States and Metrics – esxtop nomenclature



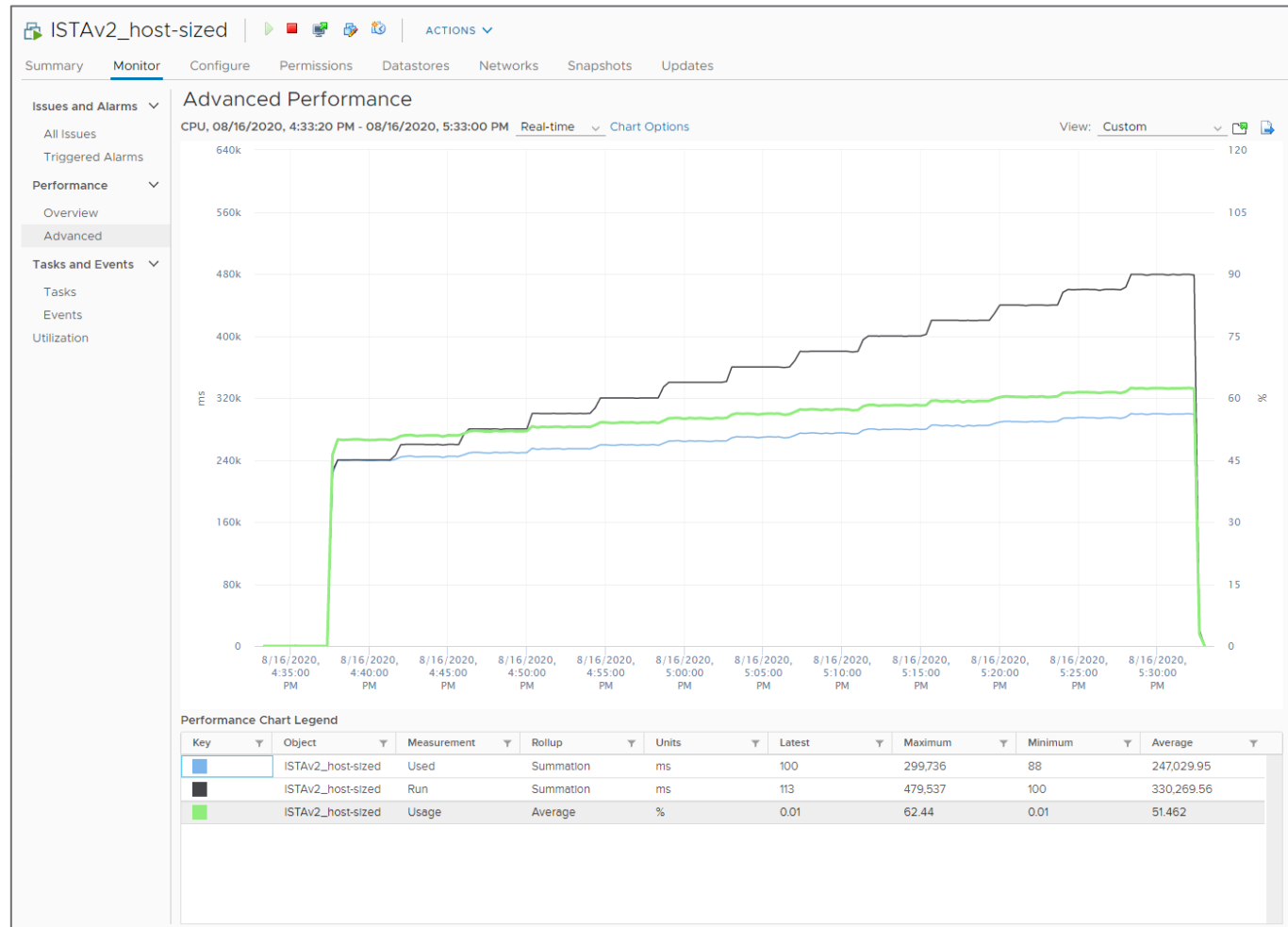
vCPU accounting

States and Metrics – esxtop nomenclature



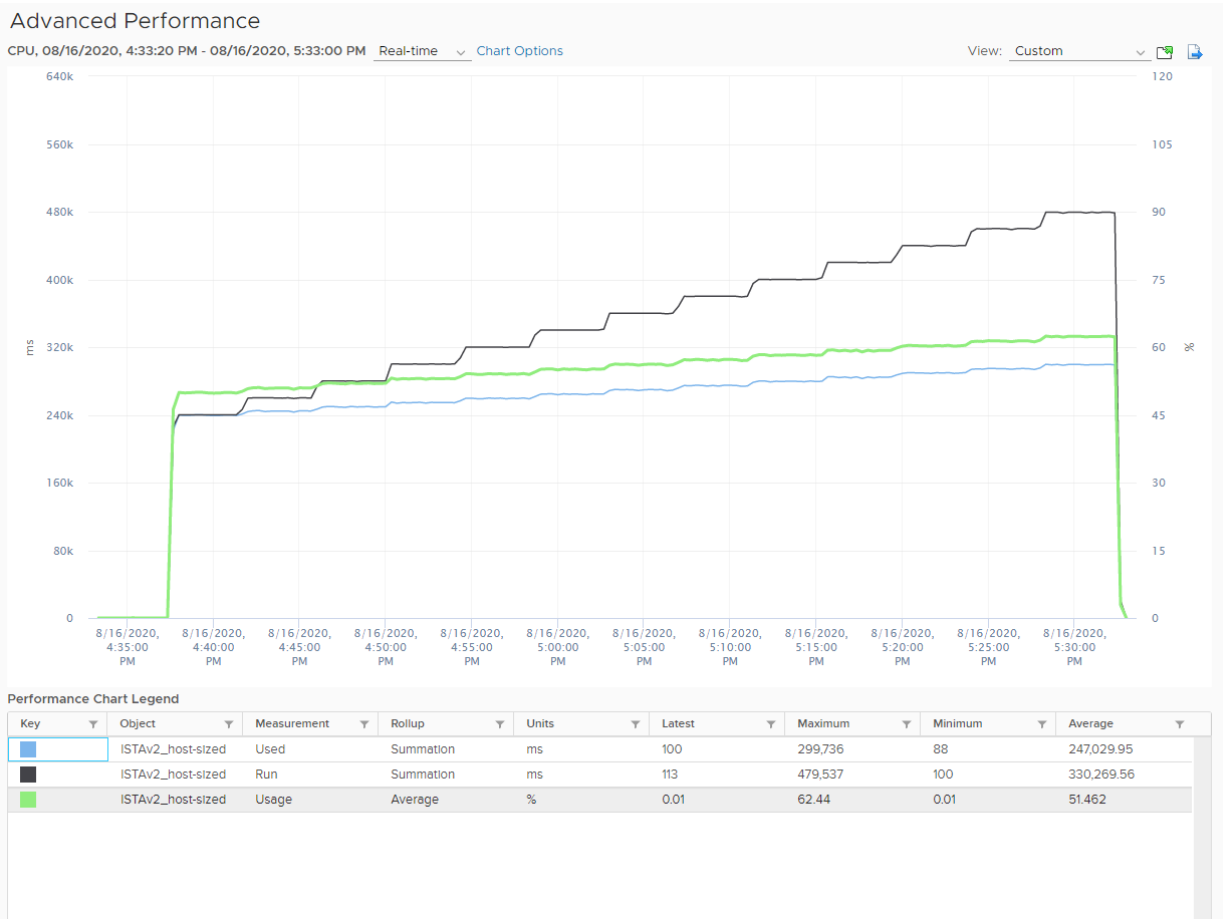
VM view

Usage only increases at a “25-degree angle”



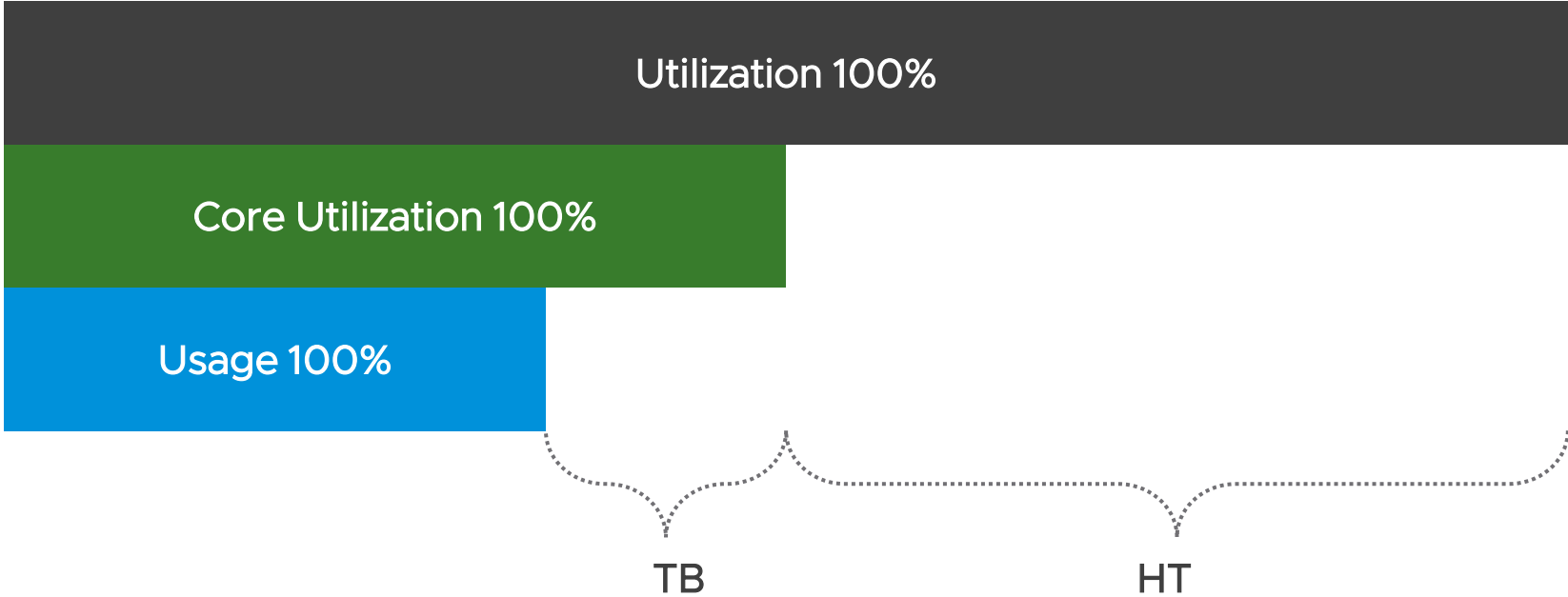
VM and Host view

Utilization / Run is the same



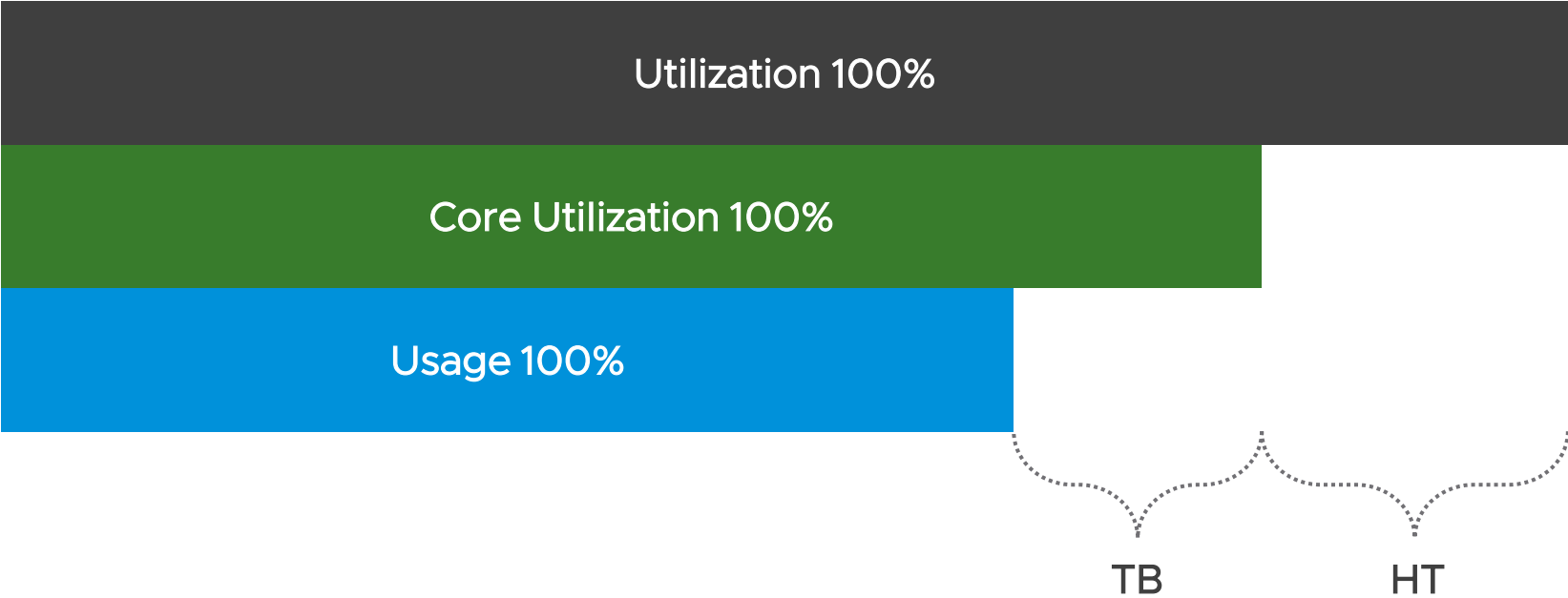
Host CPU Usage and Utilization

Looking at metrics



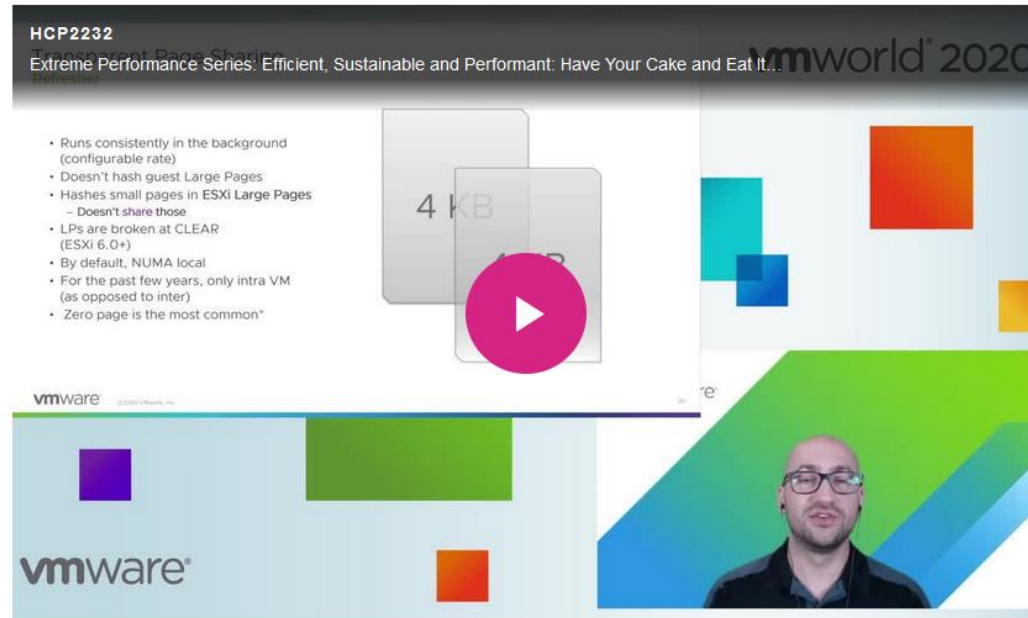
Host CPU Usage and Utilization

Looking at throughput



Host Large Pages and Transparent Page Sharing

Need to know more?



HCP2232
Transparent Page Sharing
Extreme Performance Series: Efficient, Sustainable and Performant: Have Your Cake and Eat It...

- Runs consistently in the background (configurable rate)
- Doesn't hash guest Large Pages
- Hashes small pages in ESXi Large Pages
 - Doesn't share those
- LPs are broken at CLEAR (ESXi 6.0+)
- By default, NUMA local
- For the past few years, only intra VM (as opposed to inter)
- Zero page is the most common*

4 KB

vmware

vmware

Presentation

[Download PDF](#)

Speakers



Valentin Bondzio
Senior Staff Technical Support
Engineer
VMware
[Show Bio](#)

HCP2206 - Extreme Performance Series: Efficient, Sustainable and Performant: Have Your Cake and Eat It, Too
<https://www.vmworld.com/en/video-library/video-landing.html?sessionid=1589493658174001Skyo>

Nothing's ever easy ...

Current Problems for CPU Usage / Utilization Capacity

Host Core Utilization & Host Utilization metrics

- only in Real-Time and not rolled over
 - needs to be configured manually
- not enabled by default in vROps either

Select counters for this chart:

<input type="checkbox"/> Counters	Rollups	Units	Internal Name	Stat Type	Description
<input type="checkbox"/> Ready	Summation	ms	ready	Delta	Time that the virtual machine w...
<input checked="" type="checkbox"/> Usage	Average	%	usage	Rate	CPU usage as a percentage dur...
<input type="checkbox"/> Usage in MHz	Average	MHz	usagemhz	Rate	CPU usage in megahertz during...

Timespan: **Last day** Select object for this chart:

vCenter Host and VM CPU Usage Alarms ...

Alarm Name	Object type	Defined In	Enabled
<input type="radio"/> Host CPU usage	Host	10.27.39.143	Enabled
Name: Host CPU usage			
Description: Default alarm to monitor host CPU usage			
Targets: All hosts in vCenter server 10.27.39.143			
Alarm Rules: IF Host CPU Usage is above 90% for 5 minutes THEN trigger the alarm as ! critical			
OR IF Host CPU Usage is above 75% for 5 minutes THEN trigger the alarm as ! warning			
Last modified: 02/18/2021, 5:00:25 PM			
<input type="radio"/> Virtual machine CPU usage	Virtual Machine	10.27.39.143	Enabled

Adjusting Course





Thank You

Please type your questions in the chat!