

The Architecture of the MCAO RTC for the Daniel K. Inouye Solar Telescope

*and
Matilda - a fast MVM Library for AO RTC*

Dirk Schmidt, Andrew Beard

National Solar Observatory, Boulder (Colorado), USA

ESO RTC4AO Workshop

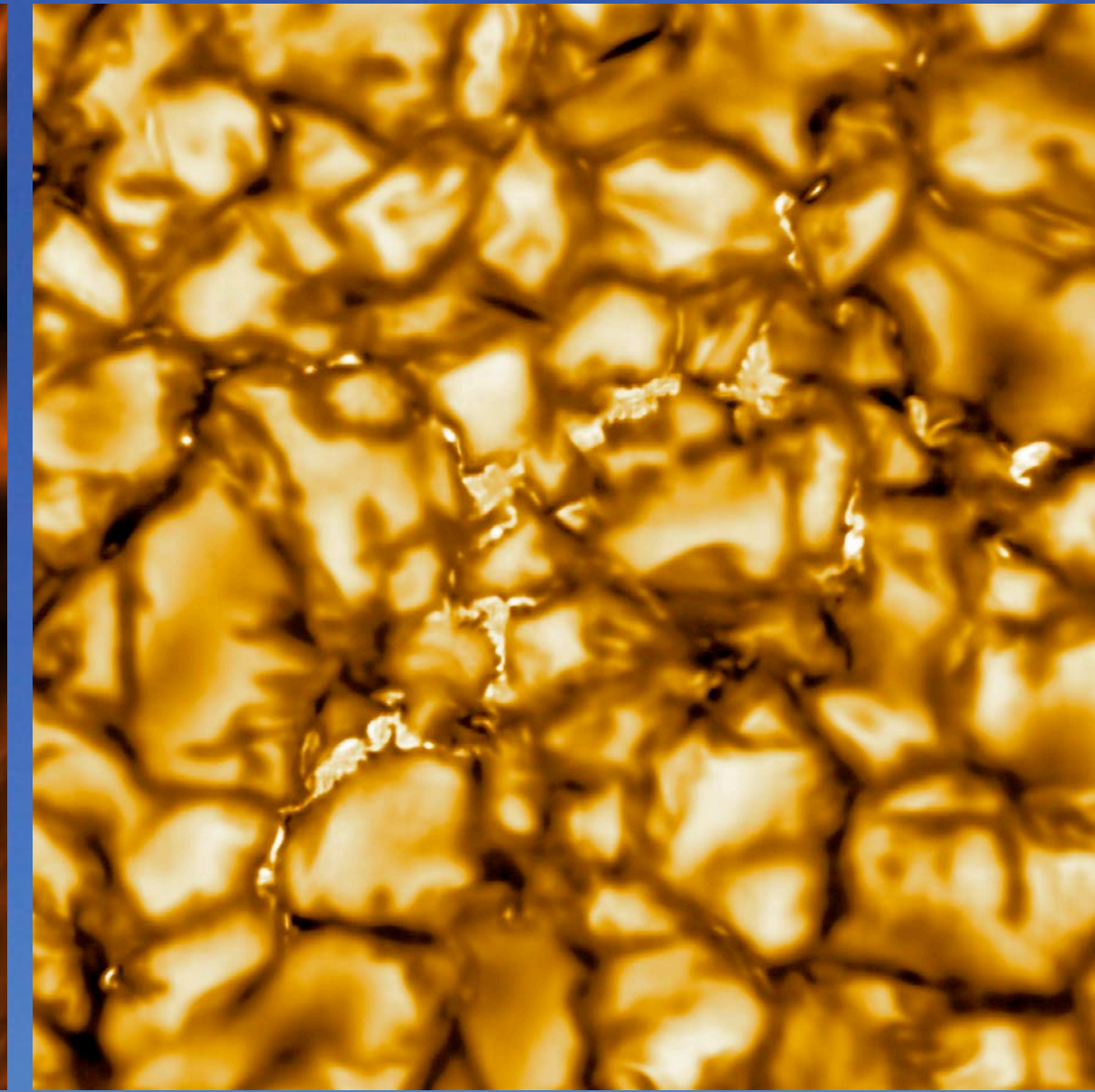
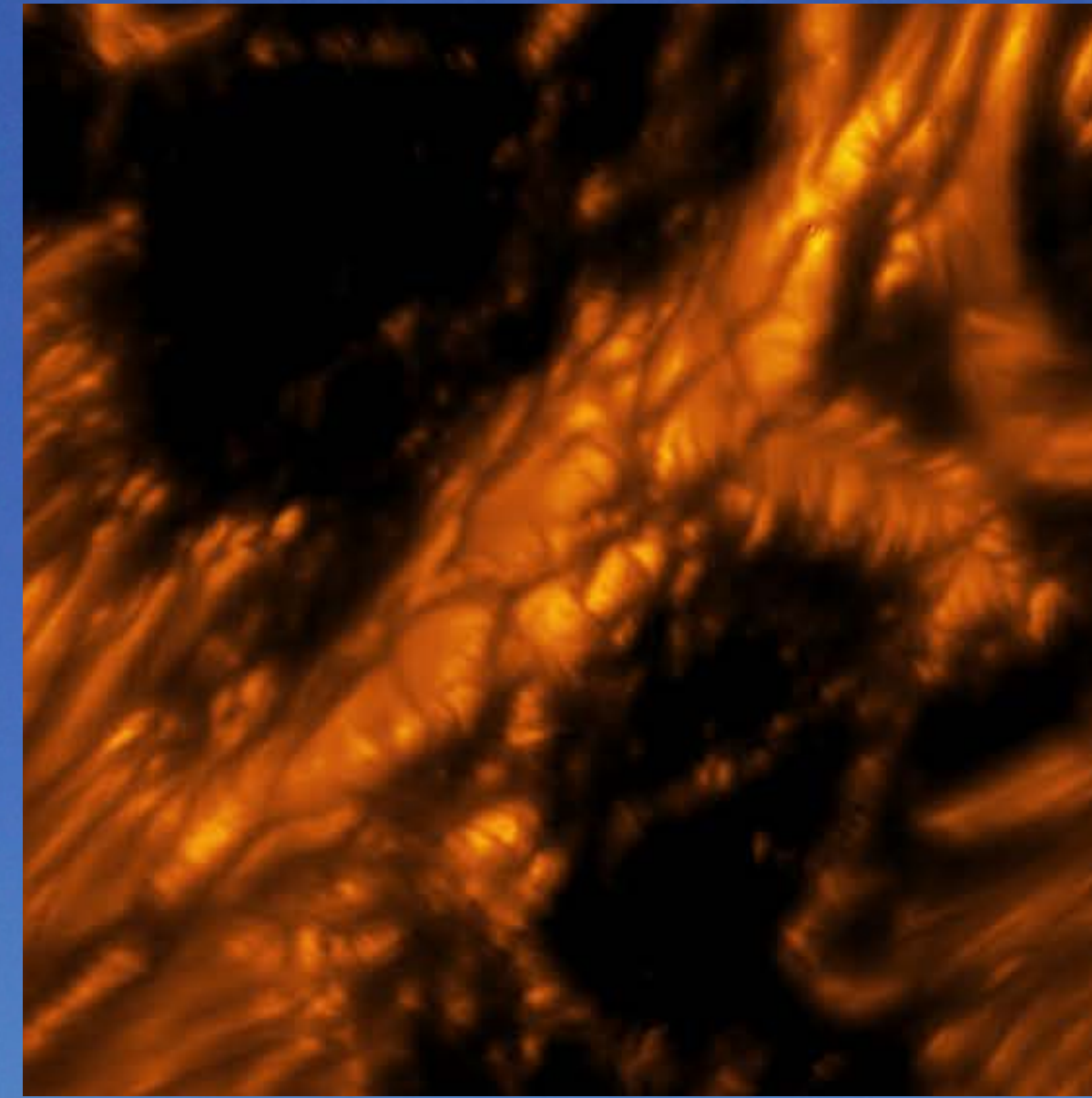
November 2023



Sponsored by the
National Science Foundation of the United States of America

The Daniel K. Inouye Solar Telescope

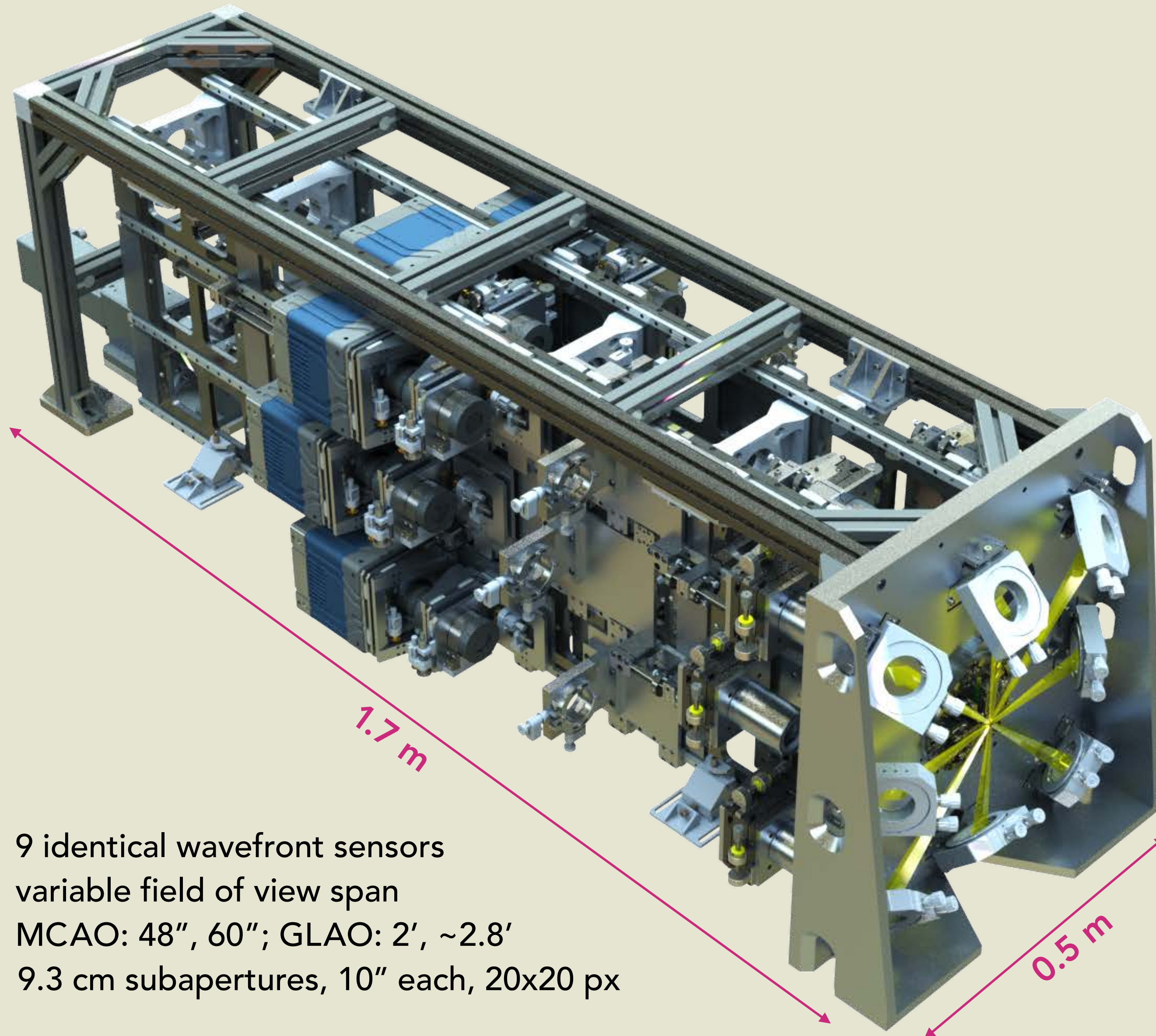
- 4-meter off-axis Gregory design (no spider)
- Science commissioning phase since 2022
- Classical Adaptive Optics with 1600 actuators
 - FPGA+CPU hybrid, 1970 fps



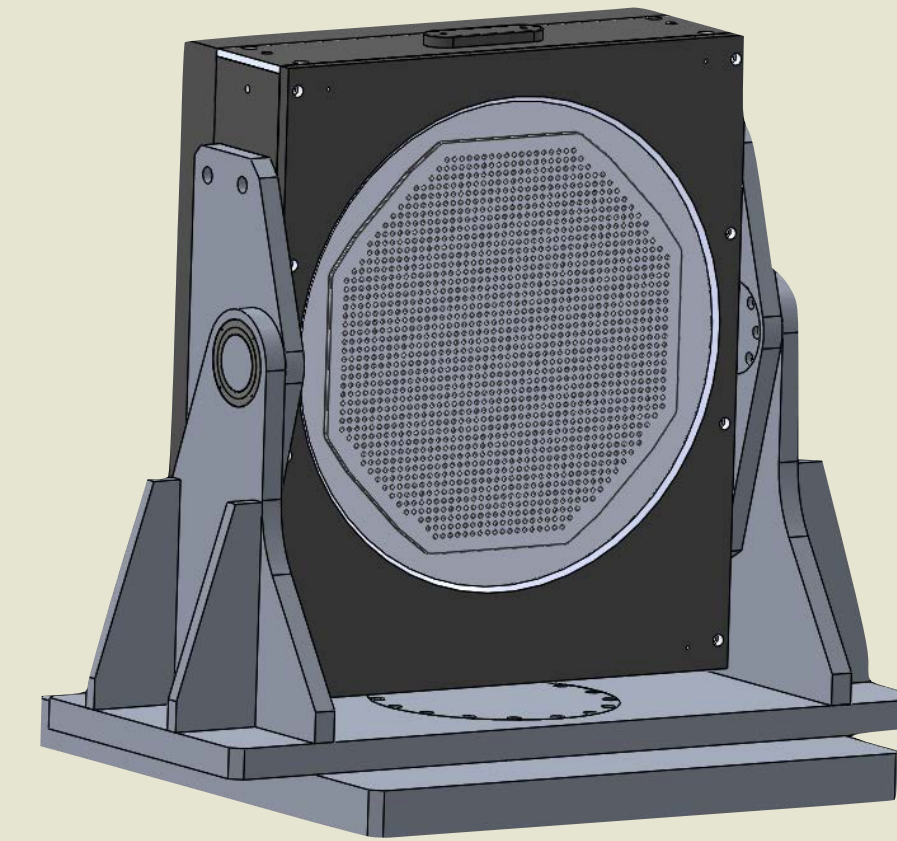
- MCAO / GLAO upgrade in fabrication
- GLAO deployment (new RTC + new wavefront sensor) by end of 2024, additional DMs later



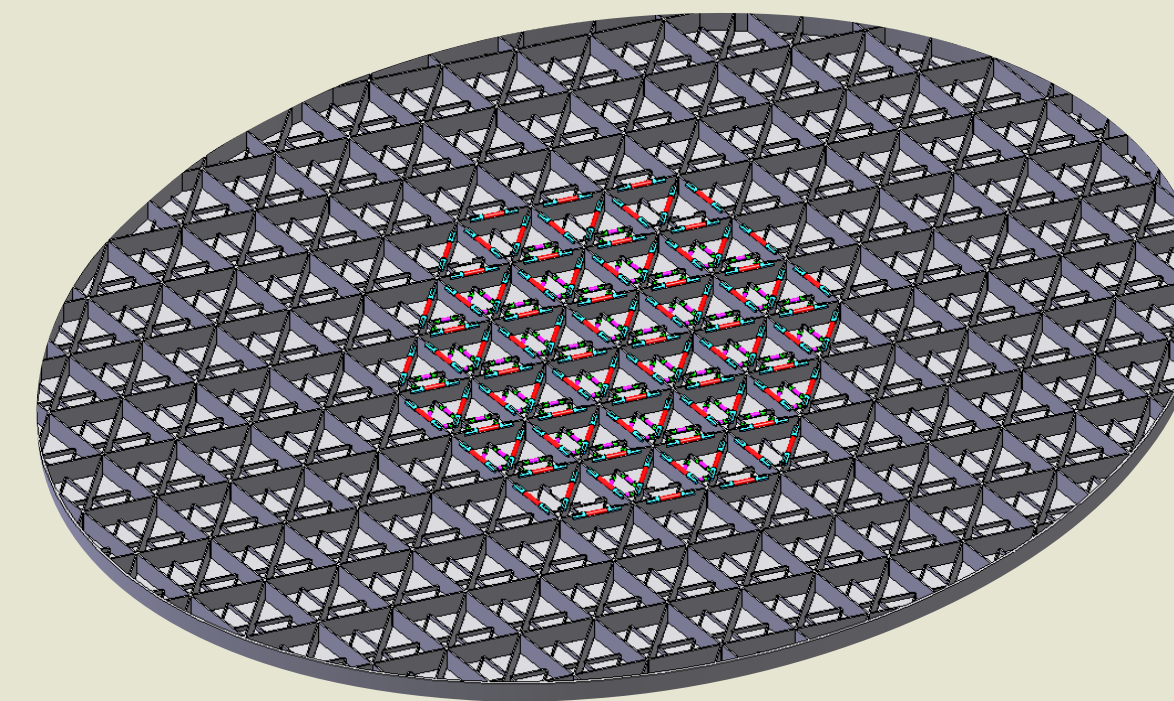
MCAO Upgrade: new WFS system, 2 additional DMs, new RTC



- 9 identical wavefront sensors
- variable field of view span
MCAO: 48", 60"; GLAO: 2', ~2.8'
- 9.3 cm subapertures, 10" each, 20x20 px



- DM_{4 km}**
- 1600 actuators
 - Final polishing



- DM_{11 km}**
- 252 actuators
 - FDR pending

11817 subapertures and 3454 actuators total



- Infiniband HDR 200 Gb/s network

- 9× Supermicro A+ 1114S-WTRT ("1of9", ..., "9of9", "The Nine")

- 1× AMD Epyc 7742 CPU (64 cores, 2.25 GHz)
- 1× Euresys Coaxlink Octo frame grabber (currently)
1× Kaya Komodo II CLHS frame grabber (in final system)
- 1× CoaXPress camera, 864×860, 8 bit, 2 kHz (currently)
1× CameraLink HS camera, 1200×1200, 12 bit, 2 kHz (in final system)
- 1× Mellanox ConnectX-6, dual-port Infiniband adapter

- 1× Gigabyte R282-Z91 ("unimatrix01")

- 2× AMD Epyc 7742 CPUs (128 cores, 2.25 GHz)
- 2× Mellanox ConnectX-6, dual-port Infiniband adapter
- 1× 10 Gb/s adapter for DM commands
- 1× 8K Nvidia graphics adapter for GUI

Real-time Control Loop Cycle

exposure

camera read out

correlations, MVM shifts to modes

camera read out

check for timing quirks

send modes to unimatrix01 via RDMA (3454)

|

Main Menu

Run AO, Check WFS Alignment..., Take Flat Images..., Ignore Subapertures..., Run AO, Calibrate AO, Check WFS Alignment..., Reference WFS, Take Dark Images..., Take Flat Images..., Calibrate System, System Service, Move Mirrors..., Move Motors..., Camera Settings..., Save Camera Snapshot..., Save Camera Burst

Run AO Menu

AD Control Mode: No correction, Image Stabilization (on-axis), Image Stabilization (full field), Classical AO, Ground-layer AO, MCAO (D = 4km), MCAO (D = 11km), MCAO (D = 4 + 11km)

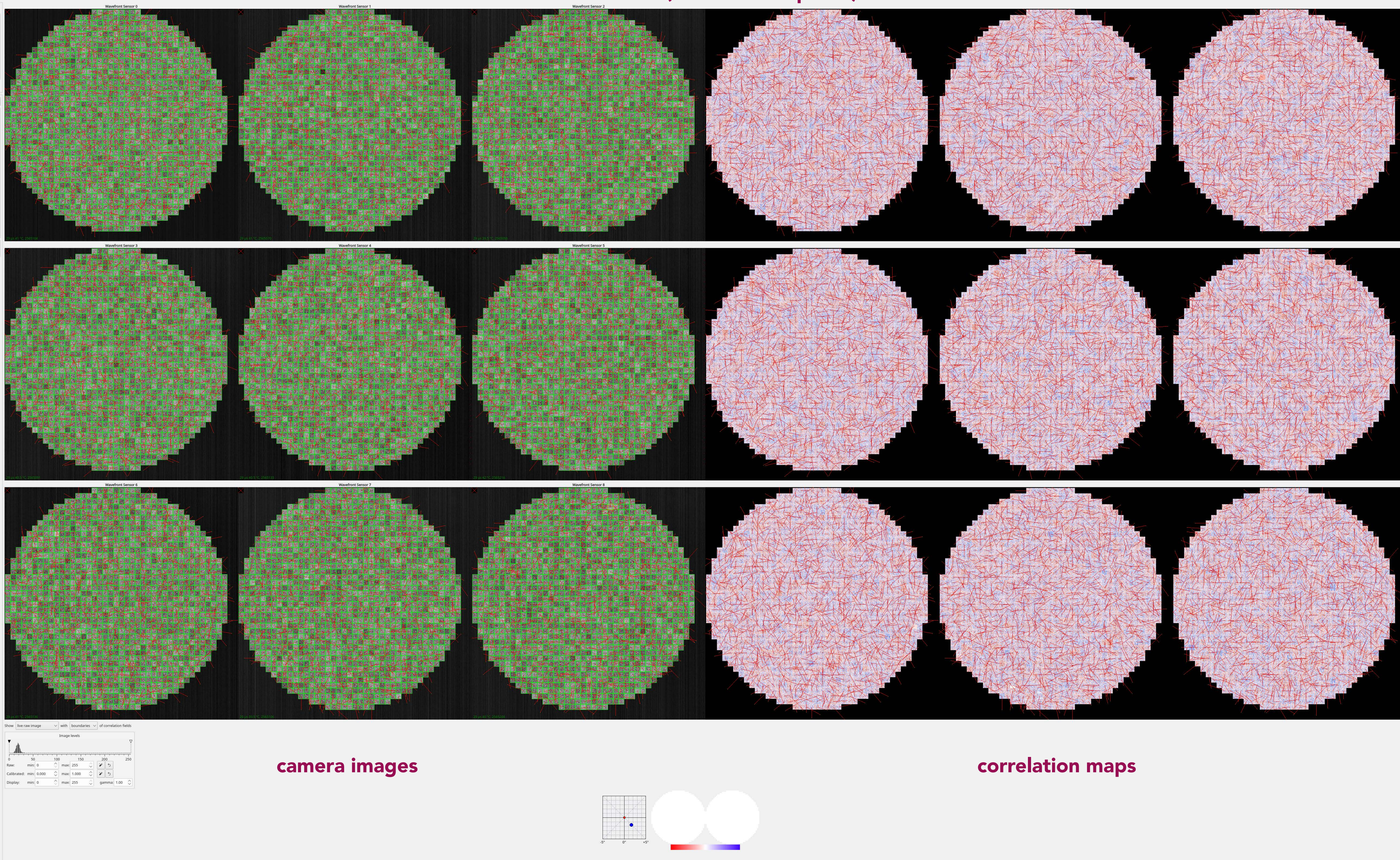
Utilities: Update Reference Image, Periodic Reference Update, Hold avg. tip-tilt position, Wavefront Sensor Settings..., Servo Parameters..., Wavefront Error Analysis..., Mode Settings..., Update Static Aberrations, Automatic Gain Optimization

Status: Mode: No correction, Seeing (r₀): 0.0 cm, Subaperture: 0.0% (#0), Contrast: 0.0000 (#0), Wavefront error after correction: mean cor shift: 3.91 μs, Active modes: D = 0 + 0 + 0 + 0, Timing: 500 μs (2000 Hz), Jitter / Latency: 16.9 μs / 40 μs, Skipped frames: 0 events, Current Buffer: 19, Reference buffer: 40%, Reference period: 0%

Display Settings: Display reduced subaperture data, Display wavefront slope vectors, Display control loop data

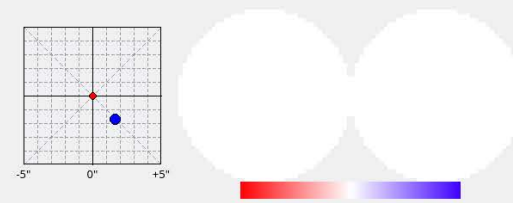
Table view				Plot view			
can	dir	subap	cor [arc]	cor [arc]	cor [arc]	act	val
0	0	1	-0.796	1.123	1.519	1	1
0	0	2	-1.881	-0.253		2	2
0	0	3	-1.839	-4.576		3	3
0	0	4	-4.487	6.427		4	4
0	0	5	-2.000	-2.000		5	5
0	0	6	-4.310	3.168		6	6
0	0	7	0.862	5.065		7	7
0	0	8	4.851	2.973		8	8
0	0	9	-4.869	-0.956		9	9
0	0	10	2.748	-6.630		10	10
0	0	11	-2.945	-1.471		11	11
0	0	12	3.114	-2.964		12	12
0	0	13	2.870	-5.287		13	13
0	0	14	4.993	3.529		14	14
0	0	15	4.818	-6.669		15	15
0	0	16	1.957	-4.123		16	16
0	0	17	2.945	-1.471		17	17
0	0	18	0.155	0.874		18	18
0	0	19	1.621	3.857		19	19
0	0	20	-6.873	-0.794		20	20
0	0	21	-6.888	6.904		21	21
0	0	22	0.449	-6.760		22	22
0	0	23	1.983	3.785		23	23
0	0	24	-1.144	-1.126		24	24
0	0	25	-1.564	0.756		25	25
0	0	26	-1.980	4.064		26	26
0	0	27	2.416	2.639		27	27
0	0	28	-4.897	-2.895		28	28
0	0	29	-1.181	-4.987		29	29
0	0	30	0.118	1.755		30	30
0	0	31	2.488	-4.187		31	31
0	0	32	-6.970	4.974		32	32
0	0	33	-0.562	-3.992		33	33
0	0	34	2.945	-1.471		34	34
0	0	35	-5.648	-6.343		35	35
0	0	36	-0.285	-2.995		36	36
0	0	37	2.817	1.795		37	37
0	0	38	2.430	-1.299		38	38
0	0	39	-4.881	2.843		39	39
0	0	40	0.221	-2.183		40	40
0	0	41	0.834	0.658		41	41
0	0	42	-3.183	5.051		42	42
0	0	43	-5.886	-2.998		43	43
0	0	44	-6.287	1.133		44	44
0	0	45	0.235	1.488		45	45
0	0	46	2.080	3.000		46	46
0	0	47	2.446	-6.843		47	47
0	0	48	-6.175	2.182		48	48
0	0	49	-1.425	2.661		49	49
0	0	50	-4.869	0.185		50	50
0	0	51	1.248	-2.185		51	51
0	0	52	0.914	-3.880		52	52
0	0	53	-6.114	-3.679		53	53
0	0	54	-6.826	6.876		54	54
0	0	55	0.991	3.887		55	55
0	0	56	0.259	3.699		56	56
0	0	57	2.814	-0.941		57	57
0	0	58	1.684	-6.263		58	58
0	0	59	-1.000	-2.000		59	59
0	0	60	-1.000	0.000		60	60
0	0	61	-5.980	1.161		61	61
0	0	62	-0.764	3.366		62	62
0	0	63	-0.864	-4.856		63	63
0	0	64	4.133	3.939		64	64
0	0	65	1.763	4.993		65	65
0	0	66	-0.014	2.800		66	66
0	0	67	-2.221	2.461		67	67
0	0	68	-1.291	0.861		68	68
0	0	69	2.129	0.179		69	69
0	0	70	3.781	-1.938		70	70
0	0	71	2.054	3.044		71	71
0	0	72	-3.935	-1.856		72	72
0	0	73	4.880	-3.008		73	73
0	0	74	3.529	-1.732		74	74
0	0	75	2.689	0.863		75	75
0	0	76	3.186	2.990		76	76
0	0	77	-0.471	4.924		77	77
0	0	78	-1.559	0.288		78	78
0	0	79	-6.800	6.800		79	79
0	0	80	0.882	-3.469		80	80
0	0	81	0.931	-9.969		81	81
0	0	82	-6.523	-6.638		82	82
0	0	83	-4.740	-1.841		83	83
0	0	84	-4.620	-3.328		84	84
0	0	85	-4.228	2.465		85	85
0	0	86	-1.000	-5.000		86	86
0	0	87	-6.598	3.969		87	87
0	0	88	-4.187	-0.214		88	88
0	0	89	-1.864	1.949		89	89
0	0	90	-0.972	2.983		90	90
0	0	91	1.073	1.962		91	91
0	0	92	-0.830	-4.655		92	92
0	0	93	-1.896	2.568		93	93
0	0	94	-1.922	-3.298		94	94
0	0	95	-1.128	3.823		95	95
0	0	96	-5.927	-0.823		96	96
0	0	97	4.854	3.838		97	97
0	0	98	0.932	-2.138		98	98
0	0	99	-1.255	-3.884		99	99
0	0	100	-4.764	-5.478		100	100
0	0	101	1.728	1.763		101	101
0	0	102	-6.800	1.000		102	102
0	0	103	-4.137	1.287		103	103
0	0	104	-3.256	3.538		104	104
0	0	105	-0.917	1.282		105	105

Image levels: Raw, Calibrated, Display, gamma: 1.00



camera images

correlation maps



RTC Software is enhanced version KAOS Evo 2



- KAOS Evo 2 originally developed for Gregor at Kiepenheuer-Institute for Solarphysics (Thomas Berkefeld & I)
- Branched off in 2013 at National Solar Observatory for Big Bear Solar Observatory
"Clear" (MCAO), AO308 II (CAO), GLAO, prominence AO
- Linux application, runs completely on x86-64 CPUs
- Complete: Control Loop and all calibrations integrated, nothing is done via external software
- GUI: designed for the operator
- Super flexible: almost everything is in config files, few performance sensitive parameters are compile time constants.
- Supports multiple SH-WFSs (cameras), each WFS can have multiple guide-regions
- Plugin-based hardware support
 - Euresys, EDT, Kaya, Active Silicon frame grabbers
 - various camera models (Mikrotron, Adimec, Photonfocus, First Light Imaging, Optronis, Dalsa)
 - Riptide Realtime DM Xinetics interface (+ obsolete in-house solutions)
- Low single-digit RMS jitter (on well tuned systems, $< 2 \mu\text{s}$ on any of The Nine)

RTC is enhanced version KAOS Evo 2

- divided into real-time back-end and GUI (Qt),
- uses both POSIX Shared Memory and threads
- “Real-time” application
 - Tuned hardware and operating system
 - Control loop threads use SCHED_FIFO with priority 94
 - Threads pinned on dedicated CPU cores
 - Threads take 100% of each assigned core and never yields the CPU (no mutexes, usleep(), etc.)
 - RMS jitter < 2 μ s (per host)

Low-latency tuning and customization for DKIST MCAO RTC

Enhancing hardware determinism and isolating CPU cores

- Consult tuning guides for HPC and low-latency applications by hardware vendors, e.g.
 - AMD *"Workload Tuning Guide for AMD EPYC 7002 Series Processor Based Servers"*
 - AMD *"Performance Tuning Guidelines for Low Latency Response on AMD EPYC-Based Servers"*
 - AMD *"Performance Tuning Guidelines for Low Latency Response on AMD EPYC™ 7002 Series Processor Based Servers"*
 - AMD *"High Performance Computing: Tuning Guide for AMD EPYC 7002 Series Processors"*
 - AMD *"HPC Tuning Guide for AMD EPYC Processors"*
- Custom BIOS from Supermicro on The Nine
 - Exposes features not found in default BIOS
- BIOS settings synchronized across servers via Redfish
 - Disabled: ASPMSupport, CorePerformanceBoost, DRAMScrubTime, HighPrecisionEventTimer, IOMMU, SMEE, SMTControl, SR-IOVSupport, SVMMode, TSME, WatchDogFunction
 - APBDIS=1, DeterminismSlider=Performance, EfficiencyModeEn=Auto, FixedSOCPstate=P0, LocalAPICMode=xAPIC

Low-latency tuning and customization for DKIST MCAO RTC

Linux tweaking, isolating CPU cores and reducing operating system jitter

- Debian Linux 11 with stock "linux-image-amd64" kernel
- unimatrix01 runs KDE desktop on 8k display (primary user interface)
- Follow "Reducing OS jitter due to per-cpu kthreads" guide
 - <https://www.kernel.org/doc/Documentation/kernel-per-CPU-kthreads.txt>
- LD_PRELOAD a custom library to prevent unwanted programs from setting their CPU affinities to isolated cores
 - overwrites sched_setaffinity() and pthread_setaffinity_np()
 - keeps "intrusive" programs which ignore preset affinities from freezing when AO RTC software is running
- All root file systems are clones from 1of9
 - Root file systems are read-only to prevent "accidental" and unrecorded changes
 - /home and /var on separate partitions, no swap partition (/var is read+write and contains symlinks for adjtime, resolv.conf, blkid.tab, etc.)

Low-latency tuning and customization for DKIST MCAO RTC

Linux tweaking, isolating CPU cores and reducing operating system jitter

```
# kernel command line
acpi_irq_nobalance apparmor=0 audit=0 clocksource=tsc cpufreq.default_governor=performance cpuidle.off=1 idle=poll
irqaffinity=0-7 isolcpus=8-63 mce=off mitigations=off nmi_watchdog=0 nohz=on nohz_full=8-63 noirqbalance nosoftlockup=0
nowatchdog processor.max_cstate=0 rcu_nocbs=8-63 rcu_nocb_poll=8-63 selinux=0 skew_tick=1 tsc=reliable

# /etc/systemd/system/cpus-offline-online.service
[Unit]
Description=Disable and re-enable CPUs 8-63
After=network.target sshd.service
[Service]
Type=oneshot
ExecStart=/bin/sh -c 'sleep 60 && chcpu --disable 8-63 && chcpu --enable 8-63'
RemainAfterExit=no
StandardOutput=journal
[Install]
WantedBy=default.target

# /etc/sysctl.d/kaos.conf
kernel.nmi_watchdog = 0
kernel.numa_balancing = 0
kernel.hung_task_timeout_secs = 0
kernel.sched_rt_runtime_us = -1
kernel.randomize_va_space = 0
vm.swappiness = 0
kernel.sched_tunable_scaling = 0
kernel.timer_migration = 0
vm.stat_interval = 10

# /etc/security/limits.d/kaos.conf
@ausers          hard    memlock    40000000
@ausers          soft    memlock    40000000
@ausers          -      rtprio     99
@ausers          -      nice       -19
```

KAOS Evo 2 design principles

- “C-style” C++
- Data layout designed (refactored) for SIMD architectures
 - “structures of arrays” where vectorization and contiguous data access is beneficial, e.g. slope, mode and actuator data
 - “arrays of structures” where not beneficial and to prevent “false sharing” of cache lines between several threads, e.g. slope, other subaperture data
 - all arrays forcefully aligned on cache line boundaries
 - not NUMA optimized (yet?), all data is local to first real-time core
- Automatic and mandated vectorization using `#pragma omp simd`
 - Analysis of vectorization reports, benchmarking and profiling (Intel ICC, Advisor, V-Tune, Compiler Explorer) to verify and improve
 - Intel ICC usually creates faster code than GCC, Clang (AOCC, Intel ICX) **even on AMD**
- Lock-free
 - no simultaneous reads and writes to the same variable by different threads
 - simultaneous writes to the same cache line are rare, too
- Thread synchronization based on busy-looping (100% CPU use)
 - only three thread synchronization barriers per loop cycle (compute shifts, modes, and actuators)
 - one atomic integer per thread, value specifies which task to execute or finished state
 - lowest overhead on AMD Epyc appears with **4 integers per cache line** despite false sharing (2 μ s for 64 threads on AMD 7742)
- Aiming for “branchless” code where it’s easy and reasonable

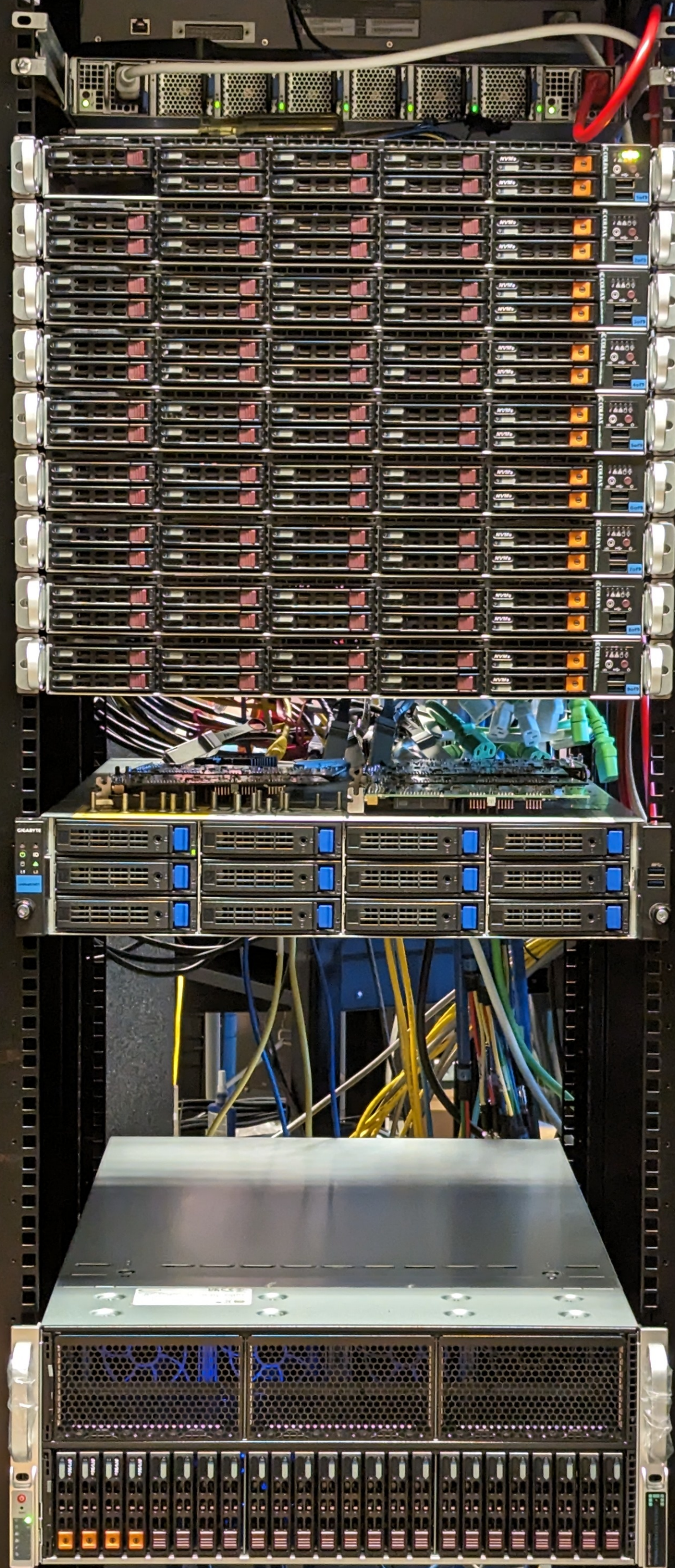
Clusterization of KAOS Evo 2

- Originally shared memory application
- Extending to cluster application
 - RDMA over Infiniband (libibverbs, librdmacm)
 - tested various write/read mechanisms to identify lowest latency approach
 - sender: RDMA_WRITE operation
 - receiver: polling raw memory (busy-looping)
 - latency: ~10 μ s to send 4000 32-bit numbers from all Nine to unimatrix01
 - RDMA layer complete, basic control loop running at 2000 Hz
 - no correlation reference updating, frequency filtering, etc. yet
 - Command layer (rpclib) almost complete
 - GUI working

Simulation-based Testing of RTC with Blur

Blur: NSO's solar AO simulator

- receives:
actuator commands
(via ethernet)
- models:
deformable mirrors,
turbulent and windy atmosphere,
imaging of extended source through atmosphere,
imaging noise
- outputs:
simulated wavefront sensor camera images
(via **CoaXPress / CameraLink simulators**, shared memory, or ethernet)
- Plugs directly into RTC instead of cameras
- Full DKIST MCAO at 10 fps on dual AMD 7742
- 1x Supermicro AS4125GS-TNRT ("chaotica")
 - 2x AMD Epyc 9654 CPUs (192 cores, 2.4 GHz)
 - 9x Kaya Chameleon II CoaXPress simulator



9 PCIe slots, 192 Zen 4 cores....

How many wavefront sensors can
KAOS Evo 2 process on chaotica?



9 PCIe slots, 192 Zen 4 cores....

How many wavefront sensors can
KAOS Evo 2 process on chaotica?

Nine!

(at 1860+ fps using 160 cores)

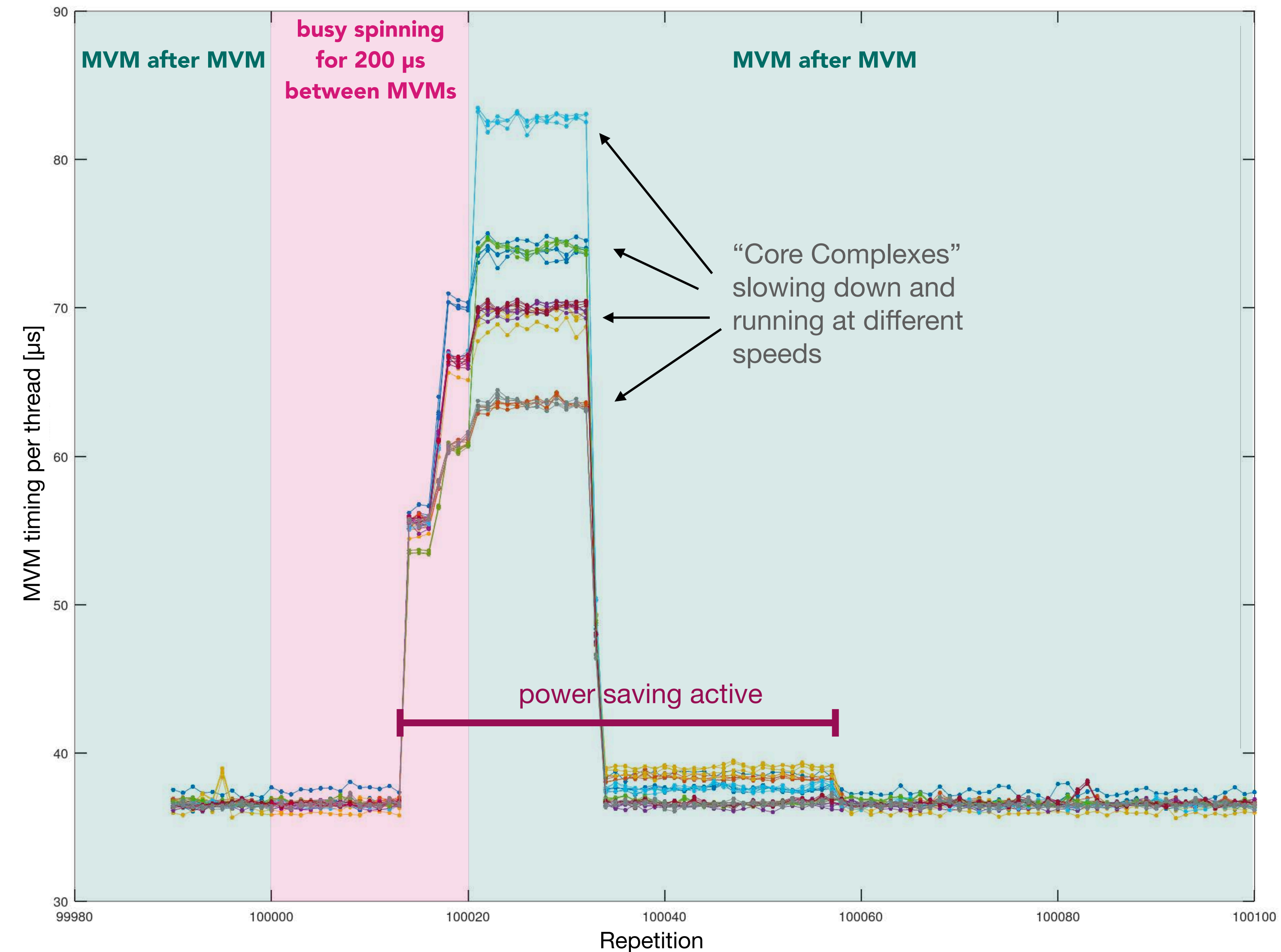
Now considering and evaluating switching gears
and run RTC on 1 or 2 computers only.

Final MVM on GPU or unimatrix01?



Throttling in AMD's 2nd Gen Epyc CPUs

MVM slows down when busy waiting for next MVM



AMD 7742 CPUs slow down

- All power savings and frequency scalings deactivated in BIOS and Linux (system tuned for low latency)
- Only seen on AMD Zen 2 CPUs (not on Intel Broadwell, Skylake, Cascade Lake, or AMD Zen 4)
- Seems to happen inside the CPU
- CPU slows down whenever deviated from full-size MVM for a short while

- “Solutions”:
 - Keep MVM going **at full size**, disregard output
 - Costs additional latency for waiting until dummy MVM has been completed
 - Avoid Zen 2 (can anybody test on Zen 1 and Zen 3?)

Matilda - a fast, low-jitter MVM library for AO RTC

Designed and optimized for repeated MVMs with constant matrix on many-core CPUs

rows x columns	MVM timing	CPUs	threads	remarks
1024 x 1024	~4 μ s	AMD Epyc 7742	64	
	~6 μ s	AMD Epyc 7742	16	
1280 x 3500	~10 μ s	AMD Epyc 9654	32	
	~13 μ s	AMD Epyc 7742	32	
2048 x 2048	~9 μ s	AMD Epyc 7742	64	
	~10 μ s	Intel Xeon Gold 6254	32	
	~25 μ s	Intel Xeon E5-2698 v4	32	
3584 x 3456 (DKIST MCAO, modes to actuators)	~15 μ s	AMD Epyc 9654	64	
	~24 μ s	AMD Epyc 7742	64	
	~27 μ s	Intel Xeon Gold 6254	32	F16C
	~67 μ s	Intel Xeon E5-2698 v4	32	
	~107 μ s	Intel Xeon Gold 6254	32	
4096 x 4096	~17 μ s	AMD Epyc 9654	128	
	~32 μ s	AMD Epyc 7742	64	
	~40 μ s	Intel Xeon Gold 6254	32	F16C
	~95 μ s	Intel Xeon E5-2698 v4	32	
	~200 μ s	Intel Xeon Gold 6254	32	
4608 x 4096	~25 μ s	AMD Epyc 7742	96	
	~38 μ s	AMD Epyc 7742	64	Intel MKL
	~105 μ s	AMD Epyc 7742	96	Intel MKL
8192 x 8192	~55 μ s	AMD Epyc 9654	128	
12800 x 3500	~24 μ s	AMD Epyc 9654	160	F16C
	~31 μ s	AMD Epyc 9654	160	
19200 x 4092	~31 μ s	AMD Epyc 9654	150	F16C
	~52 μ s	AMD Epyc 9654	150	
33280 x 8192	~105 μ s	AMD Epyc 9654	160	F16C
	~820 μ s	AMD Epyc 9654	160	

- Matilda is not a general purpose BLAS gemv() interface
MVM is usually memory bandwidth bound, but BLAS (=contiguous block of matrix data) is not bandwidth optimal on modern many-core CPUs with NUMA caches and RAM
- Matilda: First optimize data layout, then optimize the computation
Create an "mvm_plan" (memory, threads) for one particular matrix in the beginning, then execute the plan for to compute the MVM.
- MVM plan is executed on dedicated, reserved CPU cores
Matrix data cannot be evicted from caches
- Parallelized using its own thread pool designed for low latency/jitter
 - Threads are pinned to fixed CPUs and scheduling priorities can be elevated
 - Threads remain active when MVM has finished
- Explicit, hand-tuned vectorization via AVX2 and AVX-512 intrinsics
- Supports F16C 16-bit floating point format for matrix data storage to minimize cache use (MVM computed in 32-bit precision)

<https://github.com/nsomatilda/Matilda>

Matilda - a fast MVM library for AO RTC

Matrix data partitioning and layout

- Matrix data partitioned for threads
 - each thread gets a horizontal slice with multiples-of-8 rows
 - number depends on selected kernel and thread number
- Slice data transposed to column-major order
 - Copied into core-local memory pages
 - aligned on cache line boundary
 - Matrix data is not contiguous anymore, fundamental difference to BLAS.
 - High setup overhead upfront
 - More efficient memory access later

Matrix stored in row-major format

	A	B	C	D	E	F	G	H	I	J	K	L
I	0	1	2	3	4	5	6	7	8	9	10	11
II	12	13	14	15	16	17	18	19	20	21	22	23
III	24	25	26	27	28	29	30	31	32	33	34	35
IV	36	37	38	39	40	41	42	43	44	45	46	47
V	48	49	50	51	52	53	54	55	56	57	58	59
VI	60	61	62	63	64	65	66	67	68	69	70	71
VII	72	73	74	75	76	77	78	79	80	81	82	83
VIII	84	85	86	87	88	89	90	91	92	93	94	95

Matrix reordered in column-major format

	I	II	III	IV	V	VI	VII	VIII
A	0	12	24	36	48	60	72	84
B	1	13	25	37	49	61	73	85
C	2	14	26	38	50	62	74	86
D	3	15	27	39	51	63	75	87
E	4	16	28	40	52	64	76	88
F	5	17	29	41	53	65	77	89
G	6	18	30	42	54	66	78	90
H	7	19	31	43	55	67	79	91
I	8	20	32	44	56	68	80	92
J	9	21	33	45	57	69	81	93
K	10	22	34	46	58	70	82	94
L	11	23	35	47	59	71	83	95

SIMD computation on column-major buffer (SIMD length of 4 used for visualization)

	I	II	III	IV	V	VI	VII	VIII																	
I	=	0	A	1	B	2	C	3	D	4	E	5	F	6	G	7	H	8	I	9	J	10	K	11	L
II	=	12	A	13	B	14	C	15	D	16	E	17	F	18	G	19	H	20	I	21	J	22	K	23	L
III	=	24	A	25	B	26	C	27	D	28	E	29	F	30	G	31	H	32	I	33	J	34	K	35	L
IV	=	36	A	37	B	38	C	39	D	40	E	41	F	42	G	43	H	44	I	45	J	46	K	47	L
V	=	48	A	49	B	50	C	51	D	52	E	53	F	54	G	55	H	56	I	57	J	58	K	59	L
VI	=	60	A	61	B	62	C	63	D	64	E	65	F	66	G	67	H	68	I	69	J	70	K	71	L
VII	=	72	A	73	B	74	C	75	D	76	E	77	F	78	G	79	H	80	I	81	J	82	K	83	L
VIII	=	84	A	85	B	86	C	87	D	88	E	89	F	90	G	91	H	92	I	93	J	94	K	95	L

24 SIMD FMA operations only

I +=	0	A
II +=	12	A
III +=	24	A
IV +=	36	A

does not need horizontal reduction

- contiguous loading of matrix elements (very simple access pattern)
- vector elements re-used (only loaded once)
- only SIMD FMA operations used