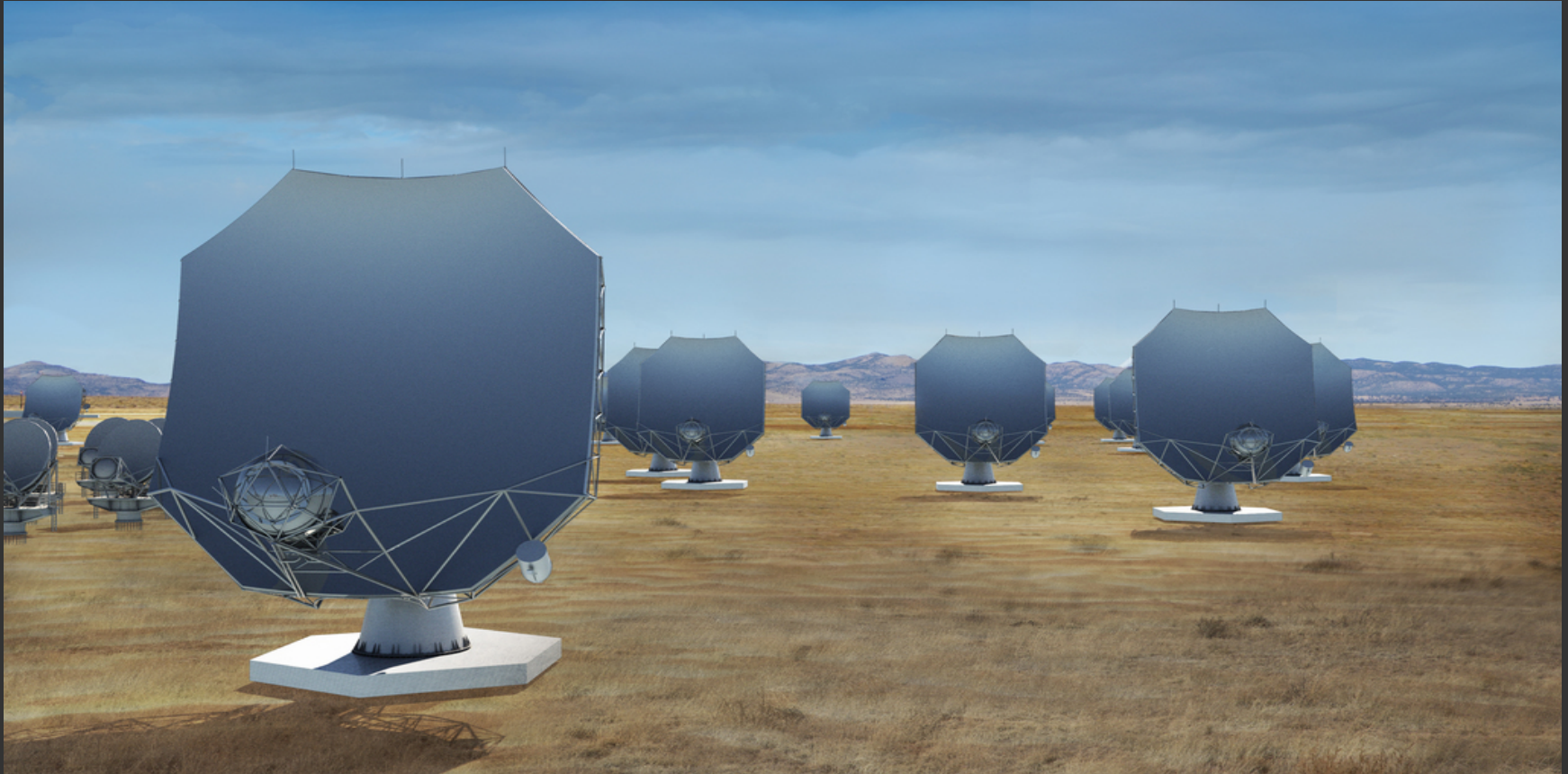# Data Processing Landscape

## PetaFLOPs, TeraBytes & Algorithms

**Wed. Lunch, Socorro, NM, May 8th 2024**



S. Bhatnagar

Algorithms R&D Group,

National Radio Astronomy Observatory, Socorro, NM, USA
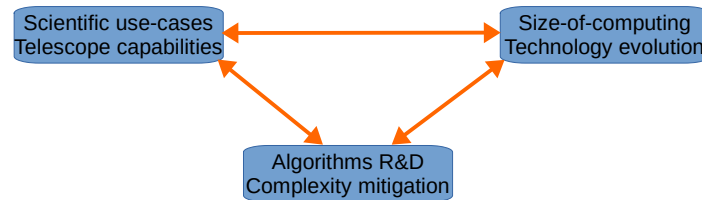
# The Algorithms R&D Group (ARDG)

- Current membership
  - Sanjay Bhatnagar (Lead) (50%)
  - Preshanth Jagannathan (Assist. Sci.) (50% ARDG, 25% CASA, 25% ngVLA CIPT)
  - Genie Hsieh (Software Eng.) (100%)
  - Felipe Madsen (100%)

- Total effort
  - ~2.5 FTE from 4 full-time staff
  - Mentoring 1 Jy PDF (Hendrik Mueller)
  - Collaborations
    - NRAO SIS Group
    - External groups/industry:
      - Kokkos(SNL/DoE), CHTC/PATh, DSA2K/CalTech, NVIDIA

# Summary

- Inherent complexity

| | |
|---|---|
| Scientific use-cases<br>Telescope capabilities | Size-of-computing<br>Technology evolution |

Algorithms R&D
Complexity mitigation

- Size of computing and the projected technology landscape

  - Scalable algorithm and software architecture

- Algorithms R&D

  - A reliable, stable software system
  - Algorithms for faster convergence, impact overall cost of computing

- Collaborations: HTC, HPC, industry groups, learn from current literature,…

  - Scaling on larger, externally managed heterogeneous clusters
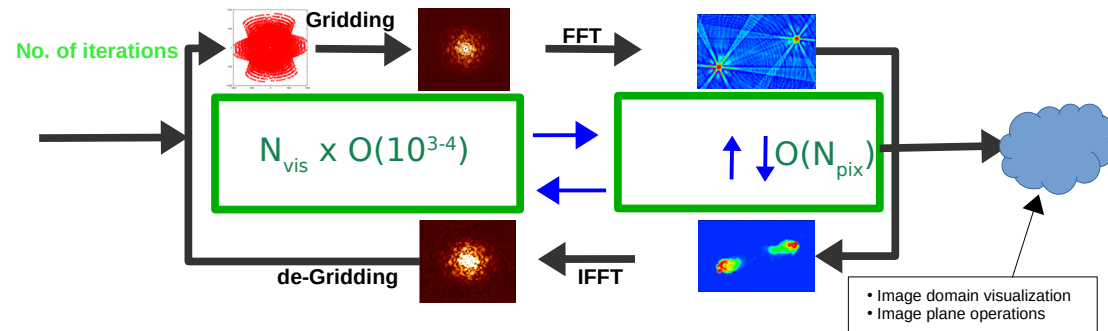  - Impact on R&D, s/w design, management,…

# System level description

- Typical data processing steps

  Imaging: $N_{vis}$ x $O(10^{3-4})$ FLOPs (Complex, SP + DP)
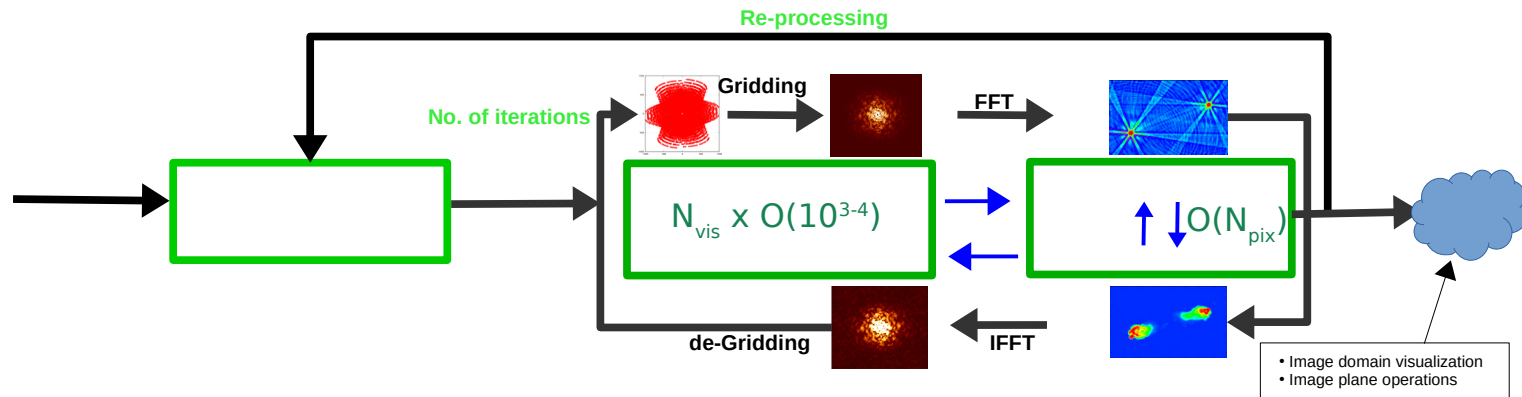  Image-plane deconvolution of the PSF : $O(N_{pix})$ FLOPs (Real-valued, SP)

# System level description

- Typical data processing steps

Imaging: $N_{vis}$ x $O(10^{3-4})$ FLOPs (Complex, SP + DP)
Image-plane deconvolution of the PSF : $O(N_{pix})$ FLOPs (Real-valued, SP)

**Re-processing**

No. of iterations | **Gridding** | **FFT**

$N_{vis}$ x $O(10^{3-4})$    $\uparrow \downarrow O(N_{pix})$

de-Gridding | **IFFT**

- Image domain visualization
- Image plane operations

# System level description

- Typical data processing steps

Imaging:                                          $N_{vis}$ x $O(10^{3-4})$ FLOPs (Complex, SP + DP)
Image-plane deconvolution of the PSF :  $O(N_{pix})$          FLOPs (Real-valued, SP)
Calibration:                                      $O(N_{vis})$          FLOPs (Complex, SP)



$$X_{ij} = \frac{V_{ij}}{V_{ij}^M}$$

$O(N_{vis})$

**Re-processing + SelfCal**

**No. of iterations**

**Gridding**

**FFT**

**Calibration**

$N_{vis}$ x $O(10^{3-4})$

$\uparrow \downarrow O(N_{pix})$

•Couples with imaging in general

**de-Gridding**

**IFFT**

- Image domain visualization
- Image plane operations

$$\forall i \rightarrow N_{ant} : initialize\left(g_i^0\right)$$
$$\forall i \rightarrow N_{ant} :$$
$$\quad d = n = 0.0$$
$$\quad \forall j \rightarrow N_{ant} :$$
$$\quad\quad d = d + g_j^{n-1} X_{ij} W_{ij}$$
$$\quad\quad n = n + \left|g_j^{n-1}\right|^2 W_{ij}$$
$$\quad g_i^n = g_i^{n-1} + \alpha \frac{d}{n}$$

**Wed. Lunch, Socorro, May 8th 2024**

# System level description

- Typical data processing steps

Imaging: $N_{vis}$ x $O(10^{3-4})$ FLOPs (Complex, SP + DP)
Image-plane deconvolution of the PSF : $O(N_{pix})$ FLOPs (Real-valued, SP)
Calibration: $O(N_{vis})$ FLOPs (Complex, SP)
Flagging: Trivial → dominant!

Re-processing + SelfCal

$f(V_{ij}, V_{ij}^M)$

$X_{ij} = \dfrac{V_{ij}}{V_{ij}^M}$

$O(N_{vis})$

Gridding

No. of iterations

FFT

**RFI mitigation**

**Calibration**

$N_{vis}$ x $O(10^{3-4})$

$\uparrow \downarrow O(N_{pix})$

de-Gridding

IFFT

- Algorithm decision not yet made
- Cost range: from trivial to more than imaging

- Couples with imaging in general

$$\forall i \to N_{ant} : initialize(g_i^0)$$
$$\forall i \to N_{ant} :$$
$$d = n = 0.0$$
$$\forall j \to N_{ant} :$$
$$d = d + g_j^{n-1} X_{ij} W_{ij}$$
$$n = n + |g_j^{n-1}|^2 W_{ij}$$
$$g_i^n = g_i^{n-1} + \alpha \frac{d}{n}$$

- *Currently* the biggest SofC driver

- Image domain visualization
- Image plane operations

# System level description

- Typical data processing steps

Imaging: $N_{vis}$ x $O(10^{3-4})$ FLOPs (Complex, SP + DP)
Image-plane deconvolution of the PSF : $O(N_{pix})$ FLOPs (Real-valued, SP)
Calibration: $O(N_{vis})$ FLOPs (Complex, SP)
Flagging: Trivial → dominant!

**Re-processing + SelfCal**

$f(V_{ij}, V_{ij}^M)$

$X_{ij} = \dfrac{V_{ij}}{V_{ij}^M}$

**Gridding**

**No. of iterations**

**FFT**

$O(N_{vis})$

| RFI mitigation | → | Calibration | → | $N_{vis}$ x $O(10^{3-4})$ | ⇄ | $\uparrow \downarrow O(N_{pix})$ | → |

**de-Gridding**     **IFFT**

- Algorithm decision not yet made
- Cost range: from trivial to more than imaging

- Couples with imaging in general

- Image domain visualization
- Image plane operations

$$\forall i \to N_{ant} : initialize\left(g_i^0\right)$$
$$\forall i \to N_{ant} :$$
$$d = n = 0.0$$
$$\forall j \to N_{ant} :$$
$$d = d + g_j^{n-1} X_{ij} W_{ij}$$
$$n = n + \left|g_j^{n-1}\right|^2 W_{ij}$$
$$g_i^n = g_i^{n-1} + \alpha \frac{d}{n}$$

- *Currently* the biggest SofC driver

HPG imaging on V100 GPU - singlethread vs. multithread

Gridding ■ Residual cycle overhead ■ Gather ■ Weights+PSF ■ Model

17280

12960

8640

4320

0

htclean (16x parallel)

"Major cycle"

"Minor cycle"

Madsen, Robnett, Rowe
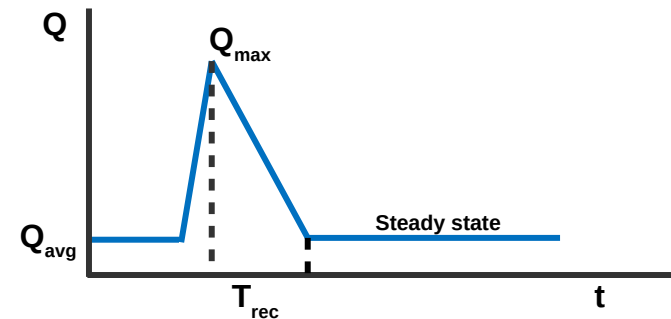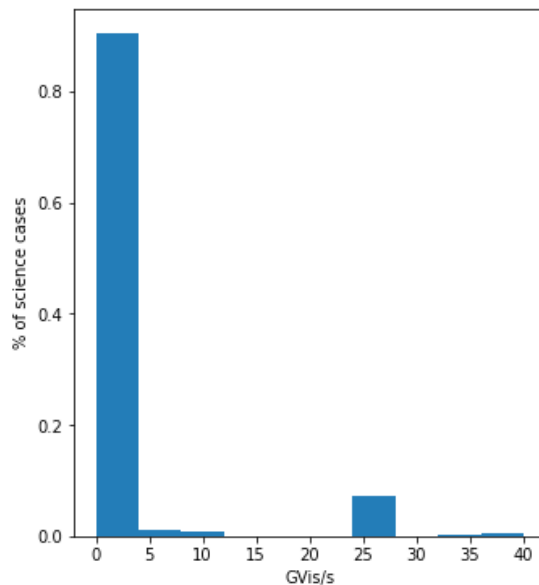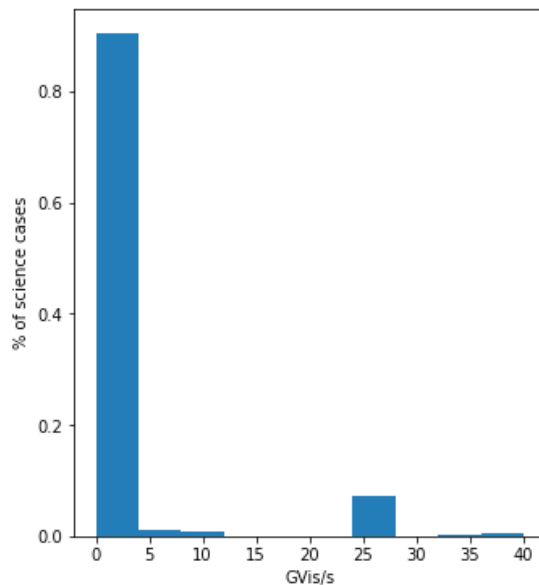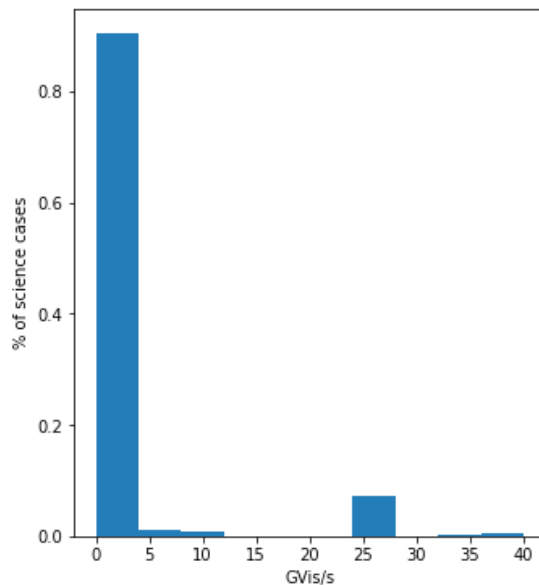
# Size of computing

- Dominated by the imaging operation    https://library.nrao.edu/public/memos/ngvla/NGVLAC_04.pdf

$$CL_{WP} = [N_{Overhead} FLOPS] \sum_{w=0}^{W_{max}-1} N_{vis}(w)[S(w=0)(\alpha w^2+1)]^2$$

$$CL_{AP} = [N_{Overhead} FLOPS] \sum_{i=0}^{N_{spw}-1} N_{vis}(\nu_i)[S(\nu_o)\frac{\nu_i}{\nu_o}]^2$$

# Size of computing

- Dominated by the imaging operation

$$CL_{WP} = [N_{Overhead} \, FLOPS] \sum_{w=0}^{W_{max}-1} N_{vis}(w) [S(w=0)(\alpha w^2 + 1)]^2$$

$$CL_{AP} = [N_{Overhead} \, FLOPS] \sum_{i=0}^{N_{spw}-1} N_{vis}(v_i) [S(v_o) \frac{v_i}{v_o}]^2$$

# Size of computing

- Dominated by the imaging operation

$$CL_{WP} = [N_{Overhead} \, FLOPS] \sum_{w=0}^{W_{max}-1} N_{vis}(w) [S(w=0)(\alpha w^2 + 1)]^2$$

$$CL_{AP} = [N_{Overhead} \, FLOPS] \sum_{i=0}^{N_{spw}-1} N_{vis}(\nu_i) [S(\nu_o) \frac{\nu_i}{\nu_o}]^2$$

$$SP_K = w_K \frac{\kappa_K \, CL_K}{\epsilon_c \, \epsilon_p} \, FLOPS/sec$$

# Size of computing

- Dominated by the imaging operation    https://library.nrao.edu/public/memos/ngvla/NGVLAC_04.pdf

$$CL_{WP} = [N_{Overhead} \, FLOPS] \sum_{w=0}^{W_{max}-1} N_{vis}(w) [S(w=0)(\alpha w^2 + 1)]^2$$
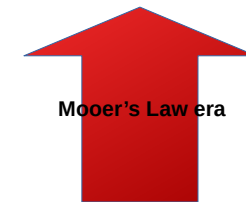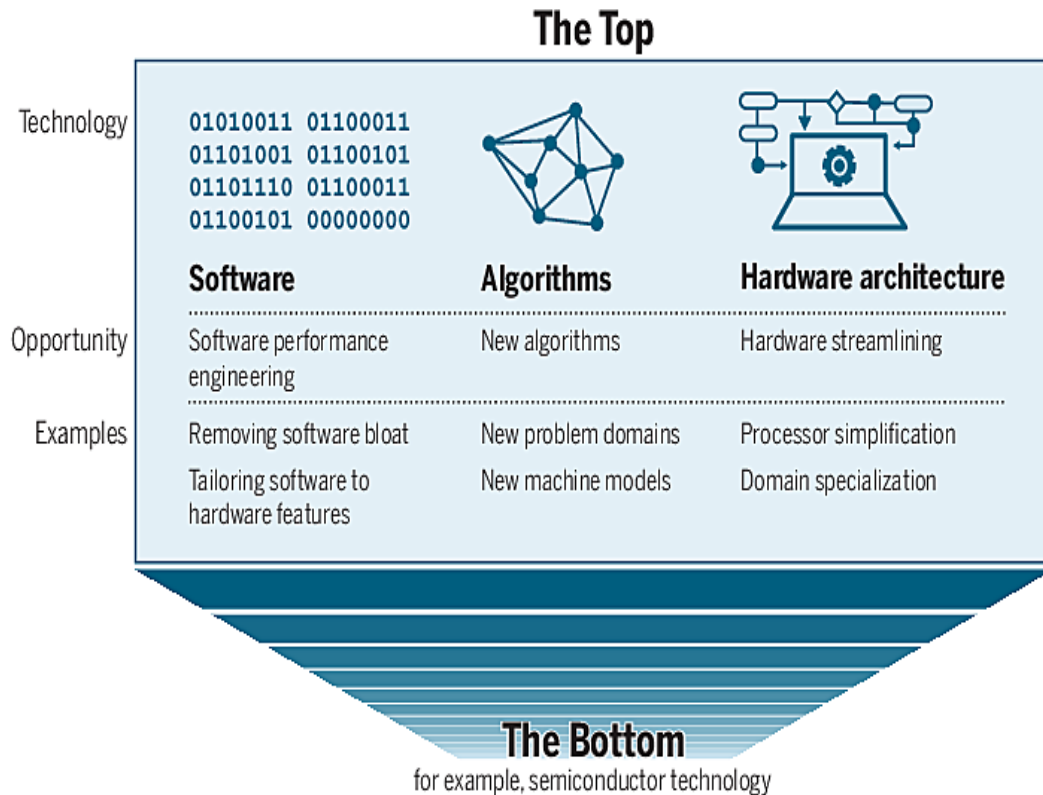
$$CL_{AP} = [N_{Overhead} \, FLOPS] \sum_{i=0}^{N_{spw}-1} N_{vis}(\nu_i) \left[ S(\nu_o) \frac{\nu_i}{\nu_o} \right]^2$$

$$SP_K = w_K \frac{\kappa_K \; CL_K}{\epsilon_c \, \epsilon_p} \; FLOPS/sec$$



- ngVLA: 50 PFLOP/sec ($T_{rec}$ ~1 day) ← ~~O(million)  CPU cores/~~few x O(1000) GPUs

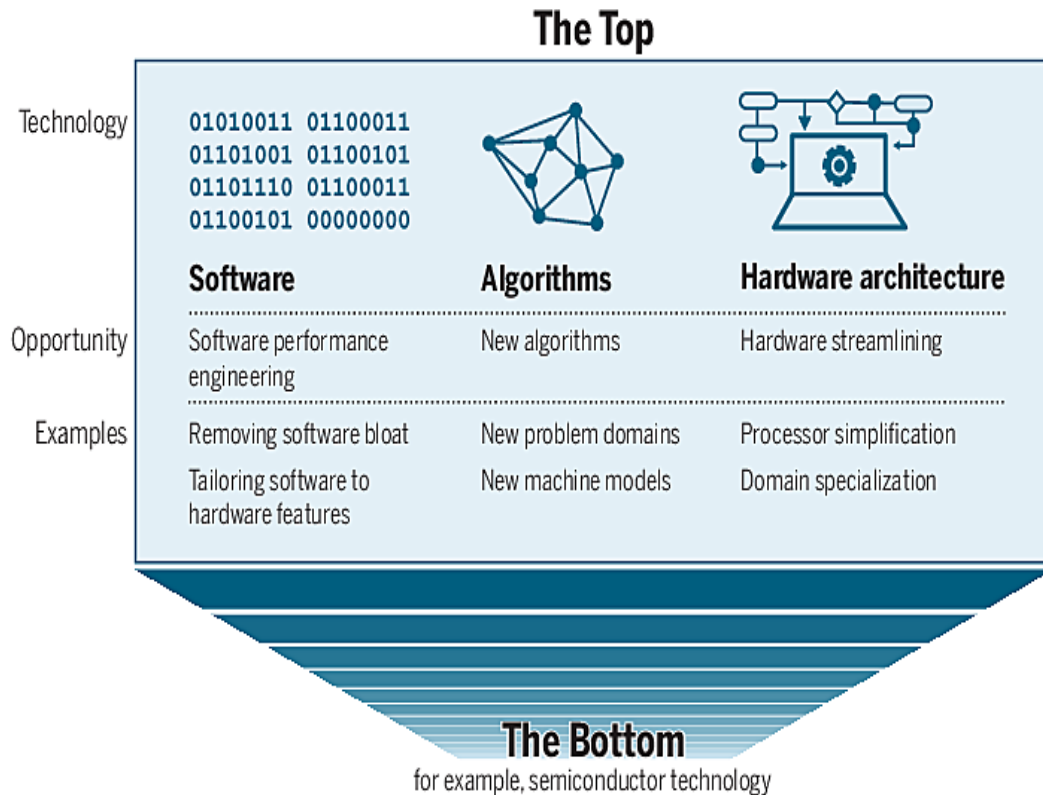- WSU:  O(100) TFLOP/sec ← Can be done today!

# Size of computing

- Dominated by the imaging operation

$$CL_{WP} = [N_{Overhead} \, FLOPS] \sum_{w=0}^{W_{max}-1} N_{vis}(w) [S(w=0)(\alpha w^2+1)]^2$$

$$CL_{AP} = [N_{Overhead} \, FLOPS] \sum_{i=0}^{N_{spw}-1} N_{vis}(\nu_i) \left[ S(\nu_o) \frac{\nu_i}{\nu_o} \right]^2$$

$$SP_K = w_K \frac{\kappa_K \, CL_K}{\epsilon_c \, \epsilon_p} \quad FLOPS/sec$$

!!

- ngVLA: 50 PFLOP/sec ($T_{rec}$ ~1 day) ← ~~O(million) CPU cores/~~few x O(1000) GPUs

- WSU: O(100) TFLOP/sec ← Can be done today!

# Computing stack: Room at the Top



**The Top**

| | Software | Algorithms | Hardware architecture |
|---|---|---|---|
| Technology | 01010011 01100011 01101001 01100101 01101110 01100011 01100101 00000000 | | |
| Opportunity | Software performance engineering | New algorithms | Hardware streamlining |
| Examples | Removing software bloat | New problem domains | Processor simplification |
| | Tailoring software to hardware features | New machine models | Domain specialization |

**The Bottom**
for example, semiconductor technology

**Performance gains after Moore's law ends.** In the post-Moore era, improvements in computing power will increasingly come from technologies at the "Top" of the computing stack, not from those at the "Bottom", reversing the historical trend.

**Mooer's Law era**

**"There is room at the Bottom"**
- Feynman (1959)

**Leiserson et al. Science (2020)**

# Computing stack: Room at the Top



**The Top**

| Technology | 01010011 01100011 01101001 01100101 01101110 01100011 01100101 00000000 | (network graph) | (hardware architecture diagram) |
|---|---|---|---|
| | **Software** | **Algorithms** | **Hardware architecture** |
| Opportunity | Software performance engineering | New algorithms | Hardware streamlining |
| Examples | Removing software bloat | New problem domains | Processor simplification |
| | Tailoring software to hardware features | New machine models | Domain specialization |

**The Bottom**
for example, semiconductor technology

**Performance gains after Moore's law ends.** In the post-Moore era, improvements in computing power will increasingly come from technologies at the "Top" of the computing stack, not from those at the "Bottom", reversing the historical trend.

**Leiserson et al. Science (2020)**

"There's plenty of room at the Top"
- Leiserson et al. (2020)

Current capacity utilization:
In single-digit percentage

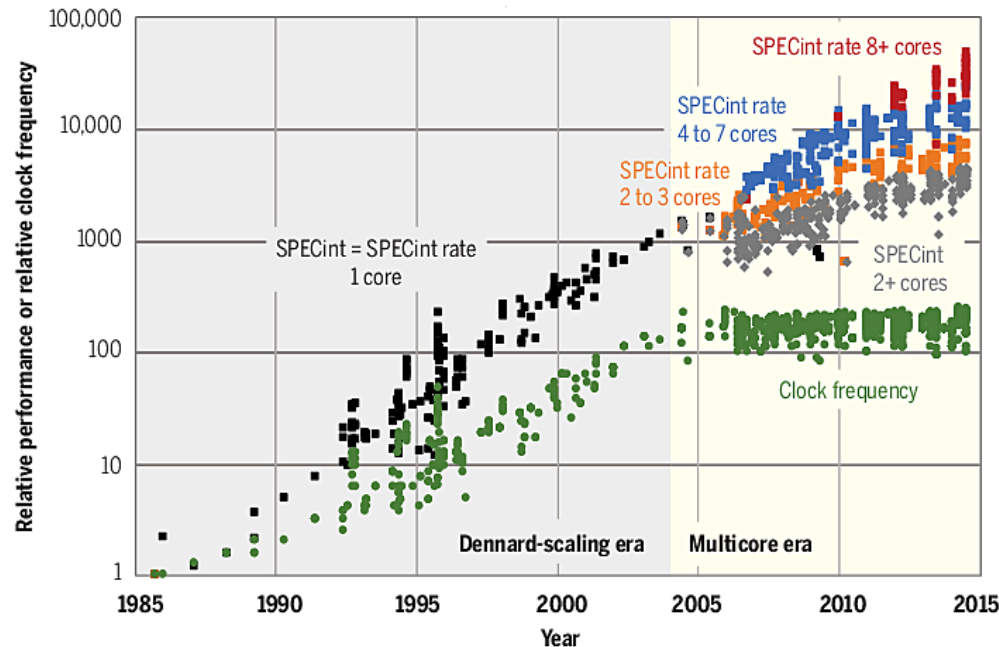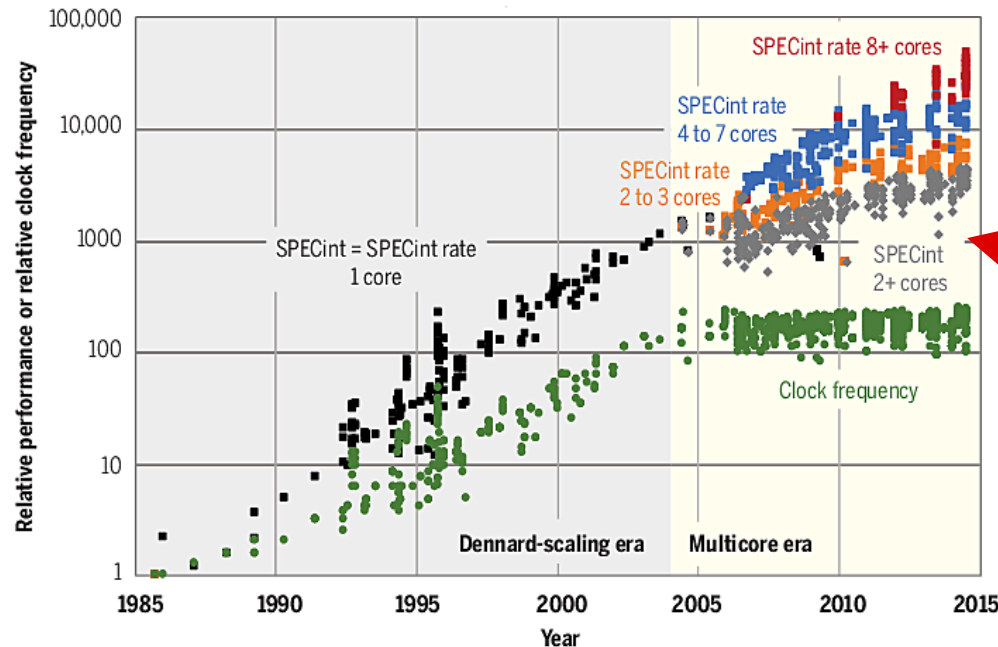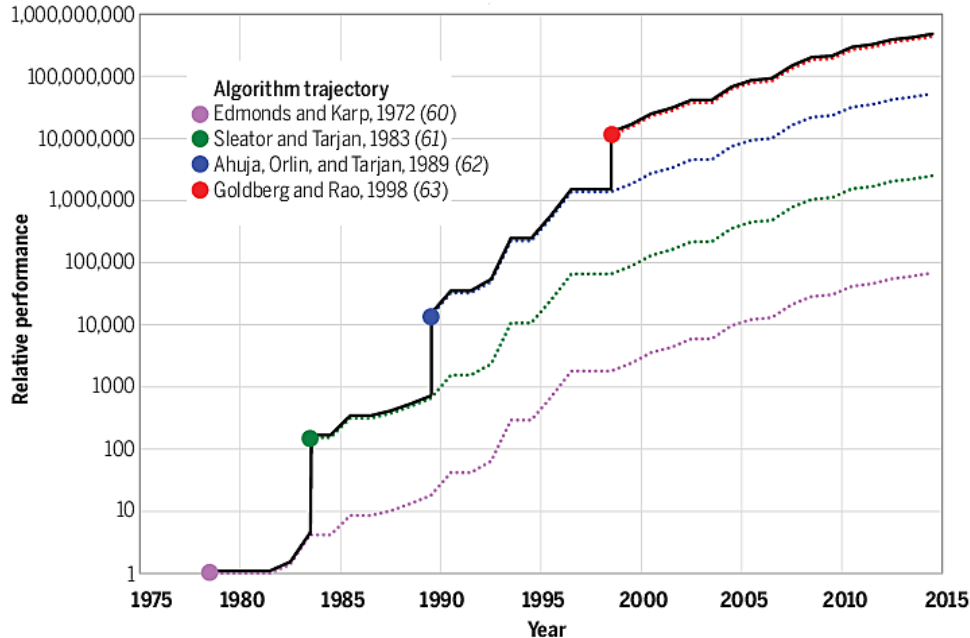Need algorithms with high
Arithmetic Intensity (FLOPs per byte)

**Mooer's Law era**

**"There is room at the Bottom"**
- Feynman (1959)

# Computing stack: Room at the Top

**The Top**

| Technology | 01010011 01100011 01101001 01100101 01101110 01100011 01100101 00000000 | (graph) | (hardware diagram) |
|---|---|---|---|
| | **Software** | **Algorithms** | **Hardware architecture** |
| Opportunity | Software performance engineering | New algorithms | Hardware streamlining |
| Examples | Removing software bloat | New problem domains | Processor simplification |
| | Tailoring software to hardware features | New machine models | Domain specialization |

**The Bottom**
for example, semiconductor technology

**Performance gains after Moore's law ends.** In the post-Moore era, improvements in computing power will increasingly come from technologies at the "Top" of the computing stack, not from those at the "Bottom", reversing the historical trend.

**Leiserson et al. Science (2020)**

"There's plenty of room at the Top"
- Leiserson et al. (2020)

Current capacity utilization:
In single-digit percentage

Need algorithms with high
Arithmetic Intensity (FLOPs per byte)

**Mooer's Law era**

**"There is room at the Bottom"**
- Feynman (1959)

# Computing stack: Multi-core era



- Moore's Law era
  - Runtime reduced by 2x if one just waited
  - Improvements were more predictable

# Computing stack: Multi-core era



- ~~Moore's Law era~~
  - ~~Runtime reduced by 2x if one just waited~~
  - ~~Improvements were more predictable~~

- Number of GP cores now is also limited by the end of Moore's-law era.

- The Top: Post Moore's-Law era:

  Massively parallel h/w of simpler cores (not GP)

  - Improvements from performance engineering, new algorithms, better silicon utilization

  - Algorithms that effectively parallelize on multiple scales of the problem

  - Specialized software

# Computing stack: Algorithms



- Historically AR&D has delivered runtime gains comparable to the Moore's Law

- Moore's Law has historically caught up...but that has now ended!

- RA algorithms have a higher FLOP per byte ratio

- RA problem: Combination of HPC (PetaFLOPs) + Big Data (TeraBytes) + 24x7 operation (High Throughput)

# Algorithm Architecture

- Stable, Scalable Architecture
  - Must scale with evolving computing needs (std VLA vs VLASS), algorithms, computing h/w & s/w (heterogeneous cluster)
  - Cast our algorithms in standard terminology: Derivative, Hessian, Update,...
  - Decompose into functionally separable components which can scale individually and together

$$V^{obs} = \boldsymbol{G^M} S \boldsymbol{F} B^M I^M + noise \qquad \chi^2 = \sum_i \left| Data_i - Model_i(P) \right|^2$$

$$P_i^{k+1} = P_i^k + [H_{ij}]^{-1} f\left(\frac{\partial \chi^2}{\partial P_i^k}\right) \quad ; \qquad [H_{ij}] = \frac{\partial^2 \chi^2}{\partial P_i^k \partial P_j^k}$$

**Model update**     **Step size**     **Derivative**

# Algorithm Architecture: Imaging

- Mathematical framework is the same for calibration and imaging
- Specialization of the components delivers various calibration and imaging algorithms



$$f\left(\frac{\partial \chi^2}{\partial I^M}\right)$$

**a.k.a. "Minor cycle"**

$$I_{k+1}^M = \left[ I_k^M + [H_{ij}]^{-1} \right] f\left(\frac{\partial \chi^2}{\partial I^M}\right)$$

**Large area networks, connected Clusters, external clusters**

**Multi-threaded CPU code**

**Residual Image a.k.a. "Major Cycle"**

# The LibRA Project: By the users, for the users

- Goals: Re-use code, re-usable library, relocatable software, ease of use

  - Derived from CASAScientific.  Now an independent code base + build system

  - Enable collaborations with RA groups and end-users  +  with other domains: HPC, HTC, Medical imaging,…

    https://github.com/ARDG-NRAO/LibRA

- Directly use the scientific layer via standalone applications
  - Deployable on external heterogeneous cluster of CPUs + GPUs

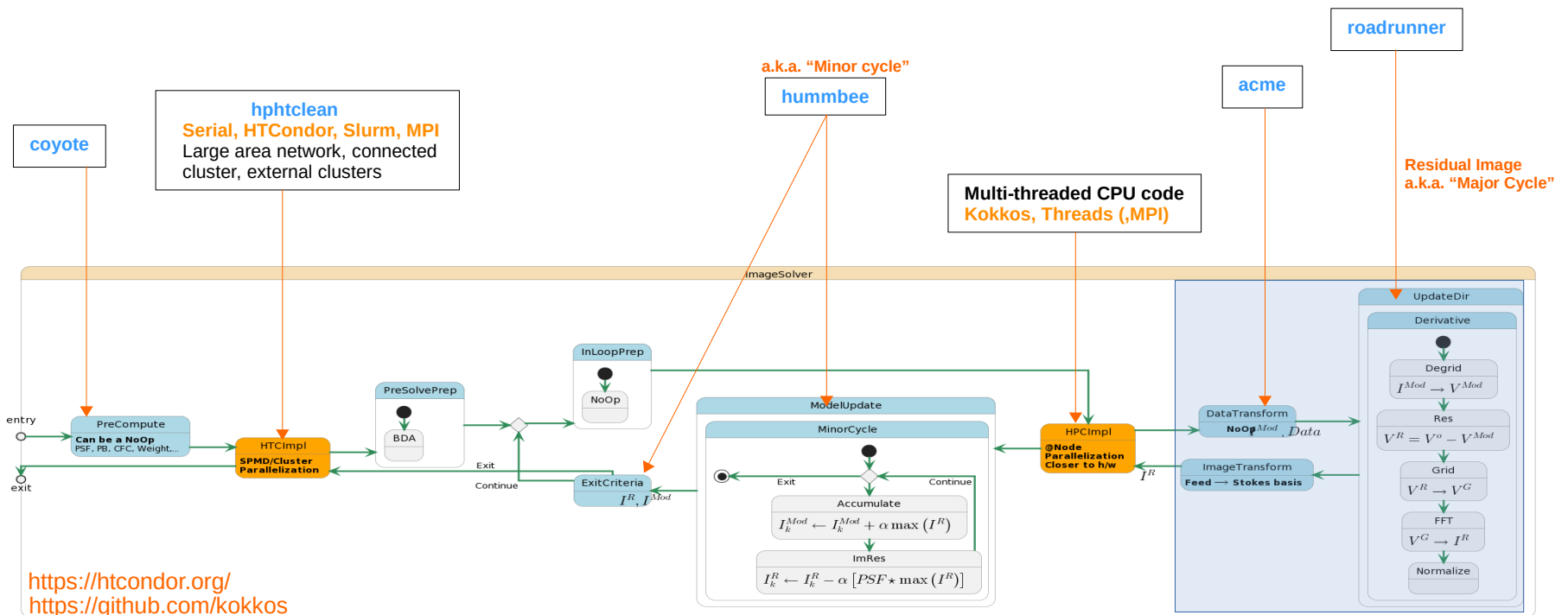- Automate chores: Modernized build system, containerized deployment, Py binding,…

# The LibRA Project: By the users, for the users

- Goals: Re-use code, re-usable library, relocatable software, ease of use

  - Derived from CASAScientific.  Now an independent code base + build system

  - Enable collaborations with RA groups and end-users  +  with other domains: HPC, HTC, Medical imaging,...

    https://github.com/ARDG-NRAO/LibRA

- Directly use the scientific layer via standalone applications
  - Deployable on external heterogeneous cluster of CPUs + GPUs

- Automate chores: Modernized build system, containerized deployment, Py binding,...



Architectural components as standalone relocatable apps

```
>roadrunner
vis                     = VLASS2.1.sb38453816.eb38509426.59047.17567765046_split.ms
imagename               = refim_oneshiftpoint.res
modelimagename          =
datacolumn              = data
sowimageext             = sumwt
complexgrid             =
imsize                  = 16384
cell                    = 0.6
stokes                  = I
reffreq                 = 3.0GHz
phasecenter             = 22:10:0.000 -00.30.0.0000 J2000
weighting               = natural
wprojplanes             = 1
gridder                 = awphpg
cfcache                 = w1.cf
mode                    = residual
wbawp                   = 1
field                   =
spw                     = 2~17
uvrange                 =
pbcor                   = 1
conjbeams               = 0
pblimit                 = 0.001
usepointing             = 0
roadrunner>
```

# Algorithm Architecture: Deployment

- Mathematical framework is the same for calibration and imaging
- Specialization of the components delivers various calibration and imaging algorithms

# High Performance Gridder (HPG)

- A gridder/de-gridder that runs on a GPUs, multi-threaded on CPUs

- Built on the Kokkos framework: Choice based on projected technology evolution
    - Implemented as a reusable independent library (ngVLA Comp. Memo #4, #5, #7)

*Tailoring software to Hardware features*

# High Performance Grider (HPG)

- A gridder/de-gridder that runs on a GPUs, multi-threaded on CPUs

- Built on the Kokkos framework: Choice based on projected technology evolution
  - Implemented as a reusable independent library (ngVLA Comp. Memo #4, #5, #7)

- Algorithm parameterized by scientific use-cases and their evolution.

$$V_{ij}^G = \left[ M_{ij}\, e^{\iota\left(\vec{\phi}_{ij}+\vec{\theta}^M\right)\cdot\Delta\vec{B}} \right] * V_{ij}^o \qquad I = FFT\left(V^G\right)$$

$$M_{ij} = CF = ATerm * WTerm * PSTerm$$

- Configurable: (Single pointing, Pointed mosaic, OTF mosaic) + Antenna pointing corrections

*Tailoring software to Hardware features*

| Operation | ATerm | PSTerm | WTerm wprojplanes | CF |
|---|---|---|---|---|
| AW-Projection | True | True False | >1 | PS*A*W A*W |
| A-Projection | True | True False | 1 | PS*A A |
| W-Projection | False | True | >1 | PS*W |
| Standard | False | True | 1 | PS |

EVLA Memo 84 (2004)
AJ, V. 154,#5 (2017)

ApJ,Vol.770, No. 2, 91 (2013)

A&A 487, 419-429 (2008)

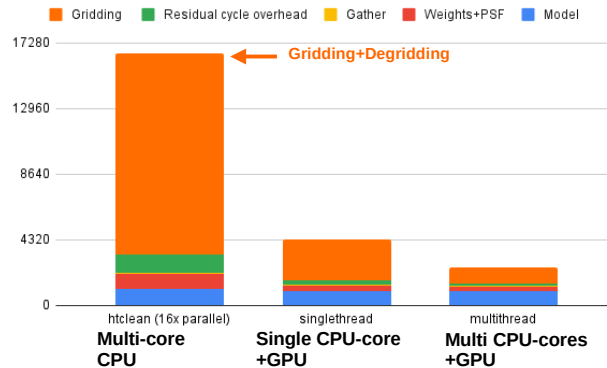# Kokkos: Performance portable eco-system (DoE)



- C++ Performance Portability Ecosystem is a production level solution for writing modern C++ applications in a hardware agnostic way.

- Part of the US Department of Energies Exascale Project – the leading effort in the US to prepare the HPC community for the next generation of super computing platforms.
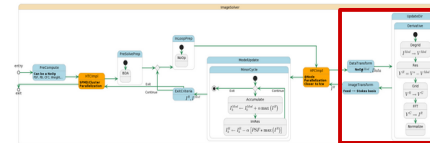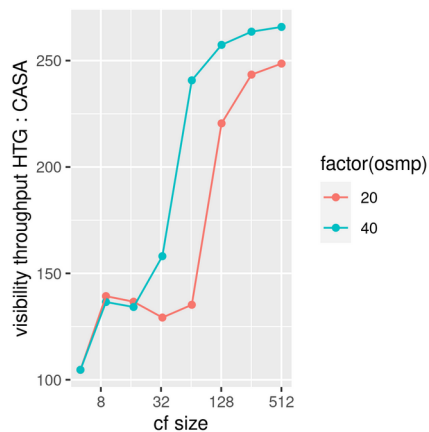
# HPG characterization

- Measured speed-up: 100 – 200x  compared to a single CPU core
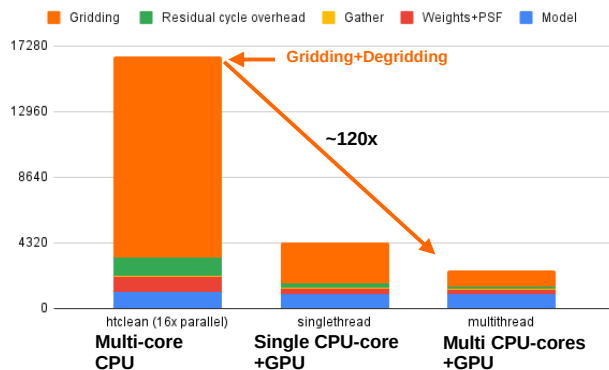


HPG imaging on V100 GPU - singlethread vs. multithread

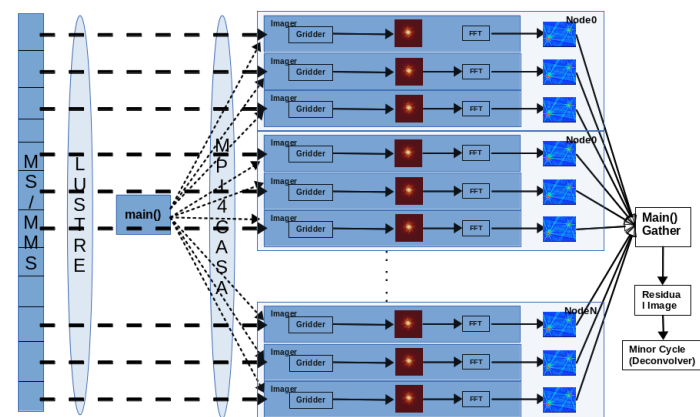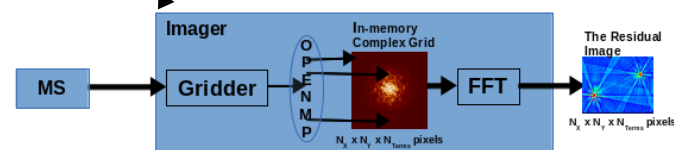Gridding — Residual cycle overhead — Gather — Weights+PSF — Model

← Gridding+Degridding

**Multi-core CPU** — htclean (16x parallel)
**Single CPU-core +GPU** — singlethread
**Multi CPU-cores +GPU** — multithread

# HPG characterization

- Measured speed-up: 100 – 200x compared to a single CPU core

HPG imaging on V100 GPU - singlethread vs. multithread

Gridding | Residual cycle overhead | Gather | Weights+PSF | Model

Gridding+Degridding

~120x

htclean (16x parallel) — **Multi-core CPU**

singlethread — **Single CPU-core +GPU**

multithread — **Multi CPU-cores +GPU**

**Complexity reduction**

$$CL_{WP} = [N_{Overhead} FLOPS] \sum_{w=0}^{W_{max}-1} N_{vis}(w) [S(w=0)(\alpha w^2 + 1)]^2$$

# Throughput measurements

- Deployed on a cluster of GPUs (100) on the PATh facility in collaboration with
  https://science.nrao.edu/enews/17.3/index.shtml#deepimaging
    - The Center for High Throughput Computing (CHTC, UW-M)
    - National Research Platform (NRP) via the OSPool, Nebraska node
    - The San Diego Super Computer Center (SDSC)
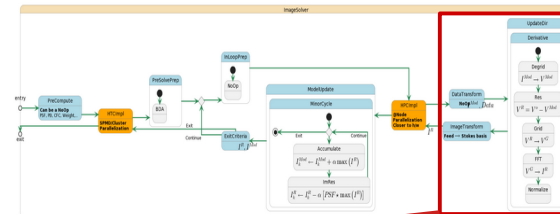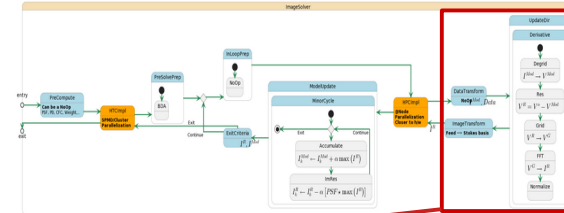    - + Multiple university computer centers across the US

# Throughput measurements

- Deployed on a cluster of GPUs (100) on the PATh facility in collaboration with
  https://science.nrao.edu/enews/17.3/index.shtml#deepimaging
  - The Center for High Throughput Computing (CHTC, UW-M)
  - National Research Platform (NRP) via the OSPool, Nebraska node
  - The San Diego Super Computer Center (SDSC)
  - + Multiple university computer centers across the US



Domain specialization

# Throughput measurements

- Deployed on a cluster of GPUs (100) on the PATh facility in collaboration with

  https://science.nrao.edu/enews/17.3/index.shtml#deepimaging

  - The Center for High Throughput Computing (CHTC, UW-M)
  - National Research Platform (NRP) via the OSPool, Nebraska node
  - The San Diego Super Computer Center (SDSC)
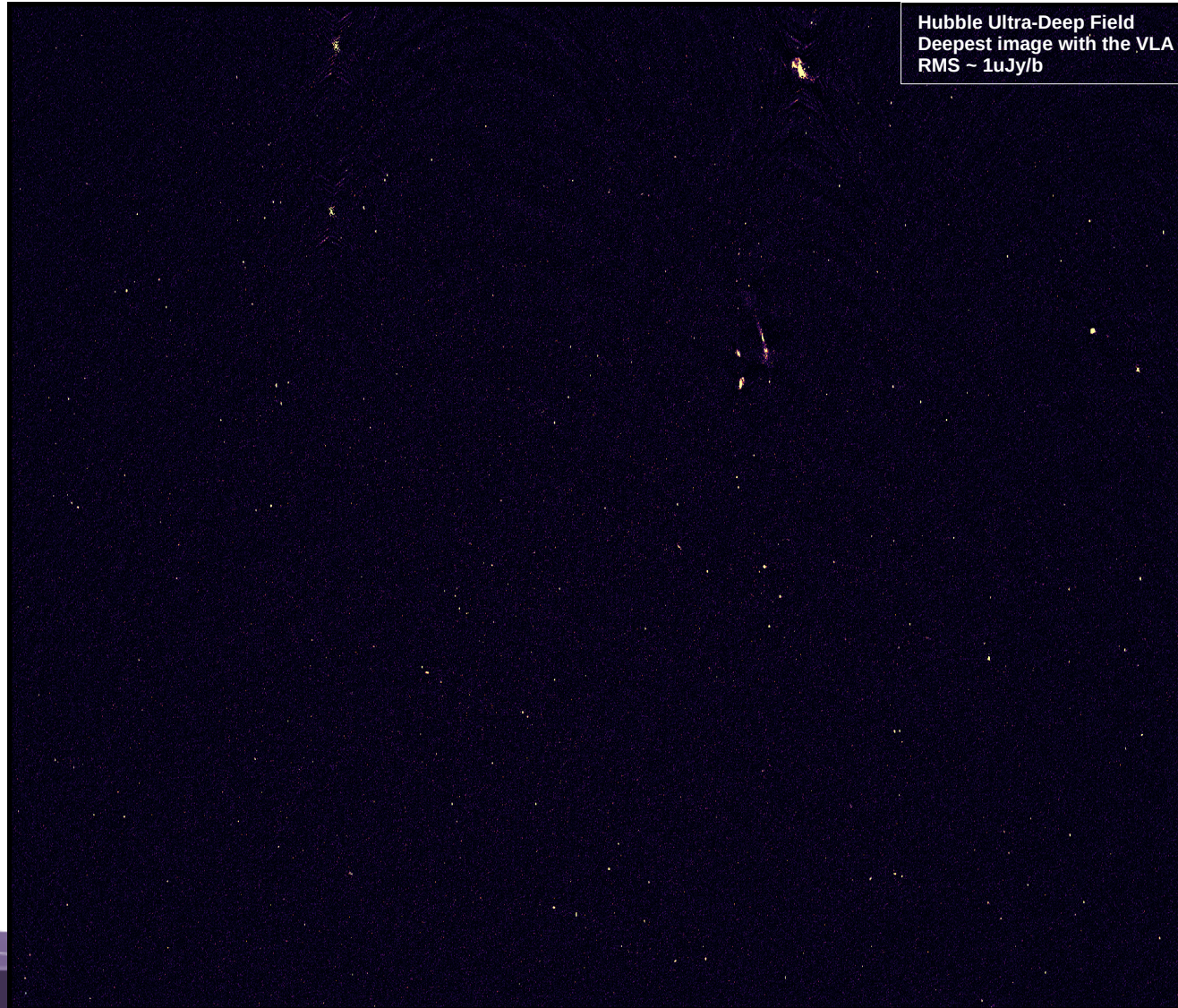  - + Multiple university computer centers across the US

Domain specialization



- Throughput: O(1 TB/hr)
- 10 iterations in ~24 hr

- Enabling tech for many unprocessed projects in the current archive:
  - Earlier attempts using CPU cores: ~14 days per cycle

- This is still a small faction of the required throughput!

**Wed. Lunch, Socorro, May 8th 2024**

# Throughput measurements

- Deployed on a cluster of GPUs (100) on the PATh facility

    https://science.nrao.edu/enews/17.3/index.shtml#deepimaging
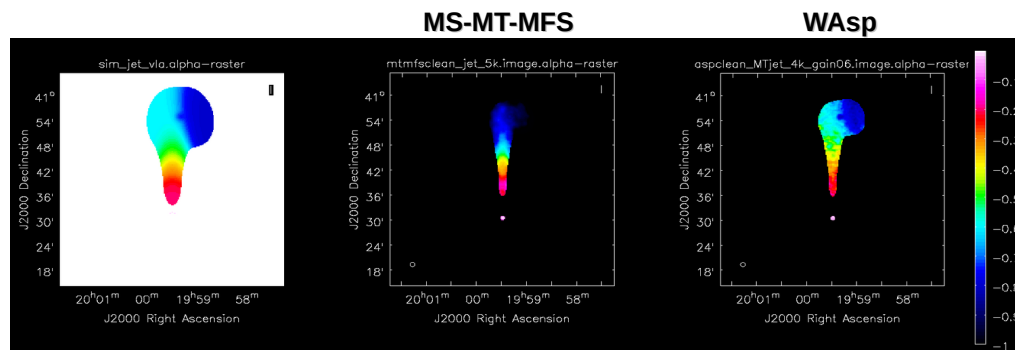


Hubble Ultra-Deep Field
Deepest image with the VLA
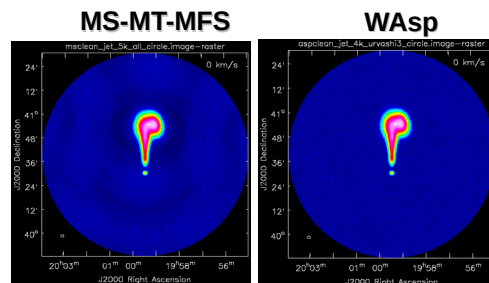RMS ~ 1uJy/b

# Image modeling (Model Update)

- Derivative calculations are most expensive → Design Model Update for faster convergence
- Scale-sensitive image reconstruction of complex emission
  - Asp-Clean : Narrow-band implementation (multi-algorithm modeling)
  - Wasp       : Wide-band Asp
  - WiS         : Wide-scale imaging (in-progress)
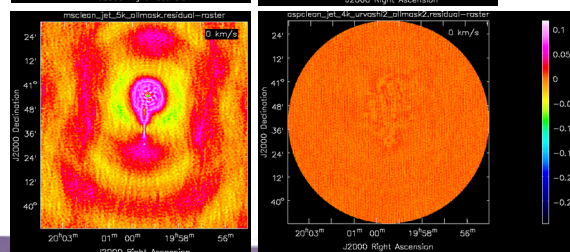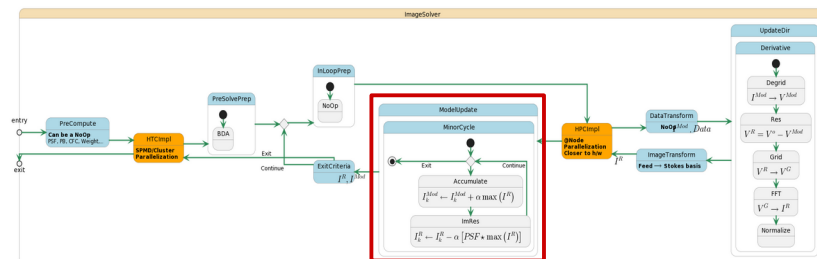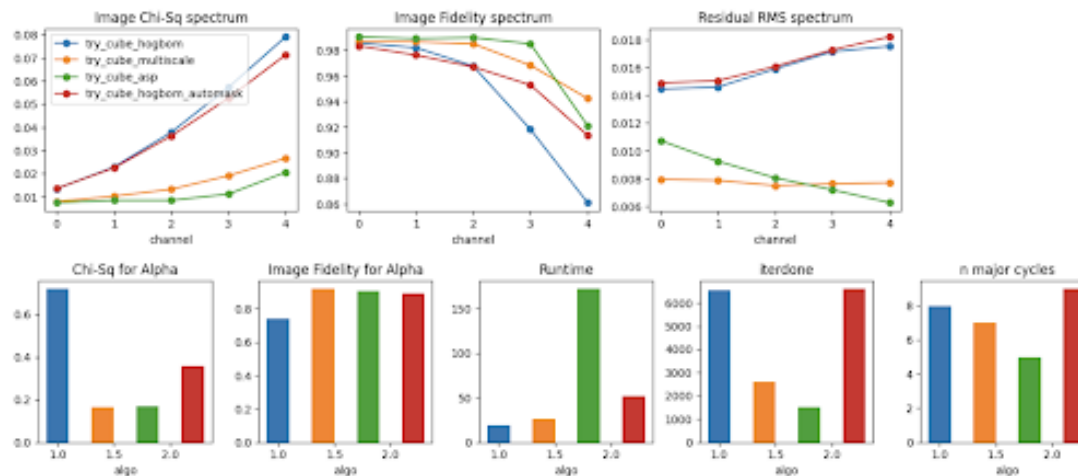
A&A, 426, 747-754, 2004

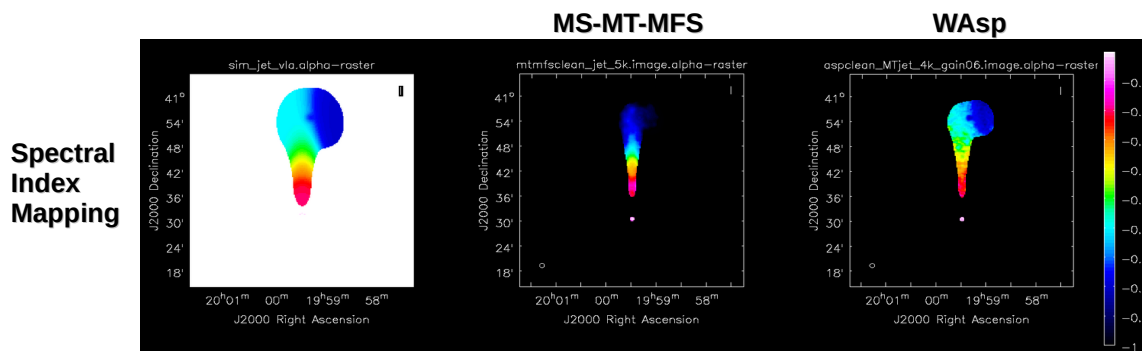# Image modeling (Model Update)

- Derivative calculations are most expensive → Design Model Update for faster convergence
- Scale-sensitive image reconstruction of complex emission
  - Asp-Clean : Narrow-band implementation (multi-algorithm modeling)
  - Wasp        : Wide-band Asp
  - WiS          : Wide-scale imaging (in-progress)

A&A, 426, 747-754, 2004



*New algorithms*

Courtesy The CASA Group

# Wide-field full-Pol. Imaging

- Wide-field full polarization mapping: The concept

$$
\begin{bmatrix} I_I^{Obs} \\ I_Q^{Obs} \\ I_U^{Obs} \\ I_V^{Obs} \end{bmatrix} = \begin{bmatrix} & & \text{ALMA DV} & \\ & & & \end{bmatrix} \cdot \begin{bmatrix} I_I^o \\ I_Q^o \\ I_U^o \\ I_V^o \end{bmatrix}
$$

**ALMA DV**
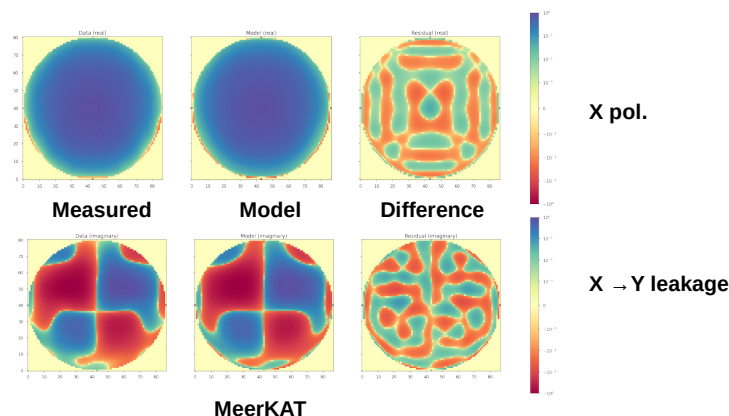
**Needs holographic measurements**

- **Diagonal**: "pure" poln. Products
- **Off-diagonal**: DD polarization leakage/mixing

Stokes-I     Stokes-V
**Before**

Stokes-I     Stokes-V
**After**

A&A 487, 419-429 (2008)

# Zernike modeling for AIP (PB)

- Build a model of the antenna aperture illumination pattern (AIP)

    - Used as input to the AW-Projection framework for wide-field full-pol. imaging. Makes the algorithmic code telescope agnostic



**ALMA DA**



X pol.

X →Y leakage

**MeerKAT**



    - Telescope agnostic tool-chain

Holography → | Plumber → *Z-coefficients* → CFCache → AWProject |

Observatory responsibility

- *Plumber* (https://github.com/ARDG-NRAO/plumber) : A general purpose package for Z-modeling of AIP, converting to PB, etc.

AJ ,163 87, 2022

# Zernike modeling for AIP (PB)

- Build a model of the antenna aperture illumination pattern (AIP)

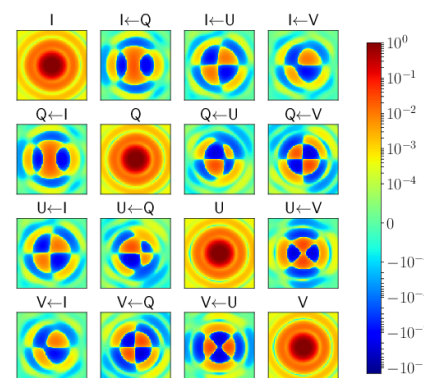  - Used as input to the AW-Projection framework for wide-field full-pol. imaging. Makes the algorithmic code telescope agnostic



**MeerKAT beam models now being used for PB correction for the MIGHTEE ~~and MALS~~ survey.**

# Conclusions thus far

- Scientific CASA code base is well designed, very re-usable, and reliable
    - ~~Is the scientific C++ code inherently as complex as imagined?~~      No.
    - ~~Is the entry-point for new developers as hard as imagined?~~        No.
        - Successful new developers/scientists in ARDG (recall – it's a 2.5-FTE group!)
          M. (Genie) Hsieh, M. Pokorny, F. Madsen, ~~S. Sekhar,~~ H. Mueller

# Conclusions thus far

- Scientific CASA code base is well designed, very re-usable, and reliable
  - ~~Is the scientific C++ code inherently as complex as imagined?~~  No.
  - ~~Is the entry-point for new developers as hard as imagined?~~  No.
    - Successful new developers/scientists in ARDG (recall – it's a 2.5-FTE group!)
      M. (Genie) Hsieh, M. Pokorny, F. Madsen, ~~S. Sekhar,~~ H. Mueller

- Minimal software stack with a robust build system reduces various costs

**Drivers**
bash, Py, Slurm, HTCondor,...

roadrunner,hummbee,acme,htclean,coyote,slurm

| libparafeed | libhpg | synthesis, {readline,ncurses,blas,lapack,gsl}-devel |
| STD C/C++ | Kokkos | RA Iterators |
| | STD C/C++ | CASACore | FFTW, WCS, CFITSIO, [HDF5, ADIOS, Boost] |
| | | STD C/C++ |

**VS**

python38 python38-devel python38-numpy, java-1.7.0-openjdk-devel python38-wheel python38-numpy python-build chrpath perl-File-Fetch python3-scipy python3-matplotlib python3-certifi python3-pytest-xvfb python3-pytest swig pkgconf-pkg-config xerces-c-devel libxml2-devel libxslt-devel sqlite-devel wcslib-devel openmpi-devel xorg-x11-server-Xvfb, redhat-lsb-core patchelf ImageMagick, Protobuf, protobuf-devel, protobuf-compiler,grpc,grpc-devel,grpc-plugin

Python shell
Tasks (Py)   Pipeline (Py)
Swig (+XML+XSLT)
Translator Layer
Py Interface C++
Tools Layer
RA ImCal Framework

| libhpg | RA Algorithms |
| Kokkos | RA Data Access/Iterators |
| STD C/C++ | CASACore |
| | STD C/C++ |

**User interface**
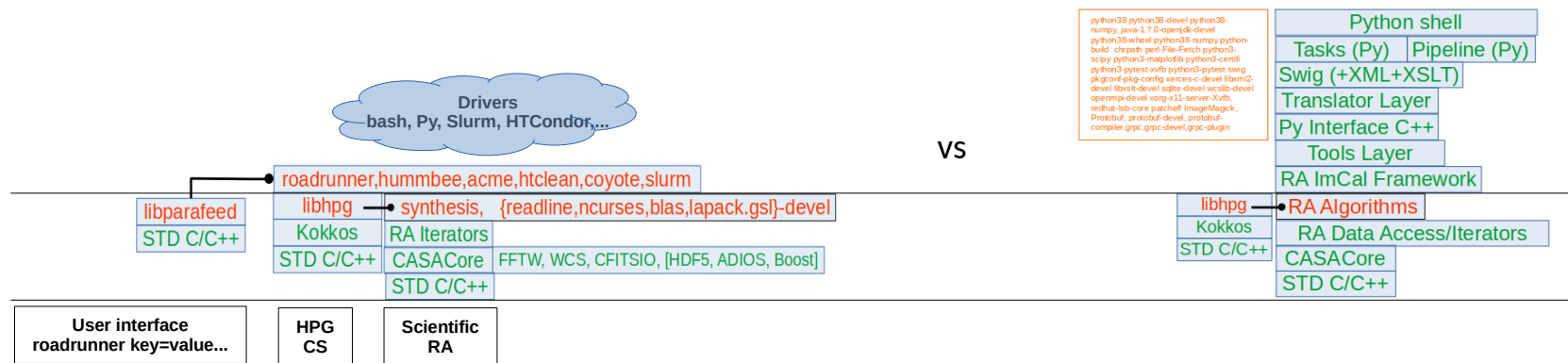roadrunner key=value...

**HPG CS**

**Scientific RA**

# Conclusions thus far

- Scientific CASA code base is well designed, very re-usable, and reliable
  - ~~Is the scientific C++ code inherently as complex as imagined?~~     No.
  - ~~Is the entry-point for new developers as hard as imagined?~~         No.
    - Successful new developers/scientists in ARDG (recall – it's a 2.5-FTE group!)
      M. (Genie) Hsieh, M. Pokorny, F. Madsen, ~~S. Sekhar,~~ H. Mueller

- Minimal software stack with a robust build system reduces various costs



- Architectural separation of functionality, development to a design, choice of technologies, and deeper understanding, keeping real use-cases (even some users) in the loop – all are important!

  - [Kokkos+libhpg]:            HPC in a h/w independent manner
  - [CASACore+libsynthesis]:  Re-use of the most advanced, highly tested RA domain scientific code-base
  - Enabling solutions:         An example of rapid deployment of scientific capability

# Work in progress

- 

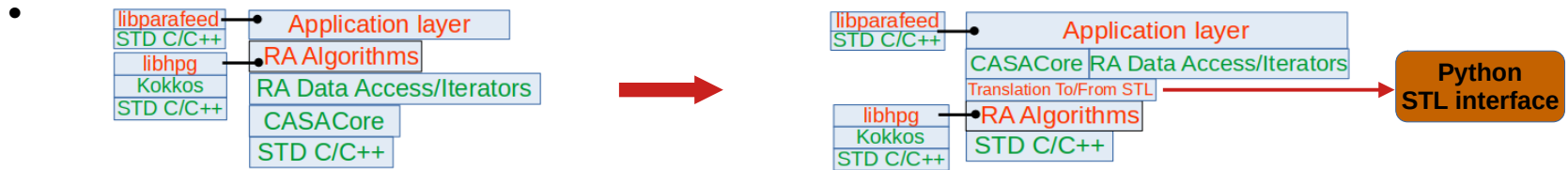**Wed. Lunch, Socorro, May 8ᵗʰ 2024**

# Work in progress

-



- Performance engineering work: NVIDIA, Kokkos/SNL,...

  - Working relationships with other groups.

    - The Kokkos group: A well established HPC R&D group that developed production code.

    - CHTC: HTC group, other communities with similar computing problem (not AI!)

    - NVIDIA (new h/w), SNL, LANL,...

# Work in progress

- 

- Performance engineering work: NVIDIA, Kokkos/SNL,…

  - Working relationships with other groups.

    - The Kokkos group: A well established HPC R&D group that developed production code.

    - CHTC: HTC group, other communities with similar computing problem (not AI!)

    - NVIDIA (new h/w), SNL, LANL,…

- Consolidate lessons from the PATh experiments

  - Getting ready for new tests + simulation + SLURM

  - Make it accessible for collaborators, other interested users

# Work in progress

- 

- Performance engineering work: NVIDIA, Kokkos/SNL,...

  - Working relationships with other groups.

    - The Kokkos group: A well established HPC R&D group that developed production code.

    - CHTC: HTC group, other communities with similar computing problem (not AI!)

    - NVIDIA (new h/w), SNL, LANL,...

- Consolidate lessons from the PATh experiments

  - Getting ready for new tests + simulation + SLURM

  - Make it accessible for collaborators, other interested users

- Algorithms R&D: Wide-scale imaging

- ngVLA Simulation, algorithm verification

  - O(100TB).  Storage is a bottleneck