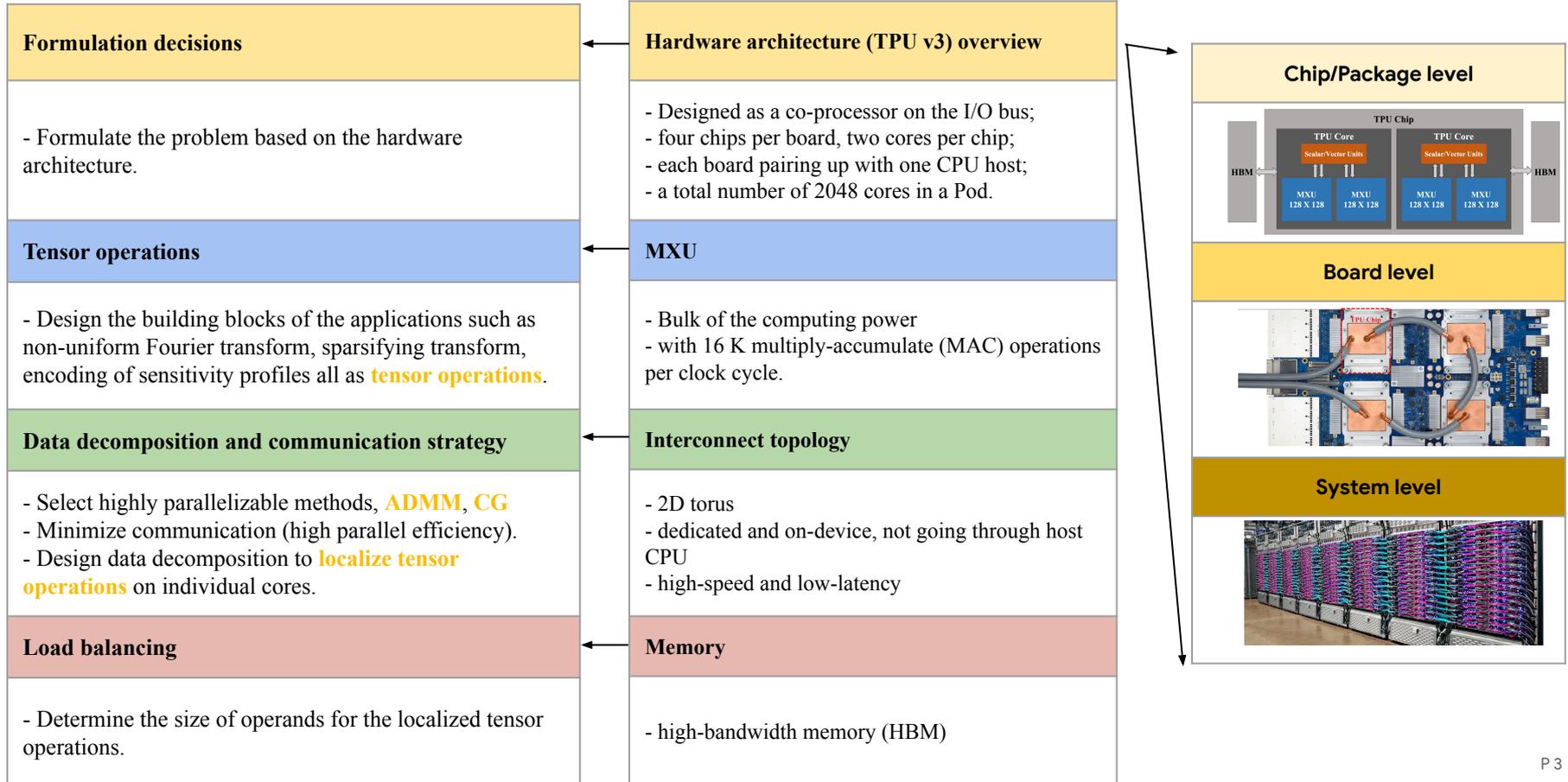


Case Studies on Accelerating Scientific Computing Applications with TPUs

Tianjian Lu and Yi-Fan Chen
June 2021

Scientific Computing on TPUs

- The recent success of deep learning has spurred the new wave of hardware accelerators.
- One such example is Google's Tensor Processing Unit (TPU).
- TPU has strength in tensor operations.
- In witnessing how scientific computing applications benefit from the advancement of hardware accelerators, it is tempting to ask whether TPU be useful for scientific computing.
- We seek an answer to this question through the case studies.



Case studies

- Fourier transform
 - Discrete Fourier transform (DFT)¹
 - Fast Fourier transform (FFT)¹
 - Nonuniform fast Fourier transform (NUFFT)²
- Linear system solver
 - Conjugate gradient (CG) method³
- Numerical optimization
 - Alternating direction method of multipliers (ADMM)³
- The applications in medical imaging³

1. Lu, Tianjian, Yi-Fan Chen, Blake Hechtman, Tao Wang, and John Anderson. "Large-scale discrete Fourier transform on TPUs." *IEEE Access* (2021).
2. Lu, Tianjian, Thibault Marin, Yue Zhuo, Yi-Fan Chen, and Chao Ma. "Nonuniform Fast Fourier Transform on Tpus." In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pp. 783-787. IEEE, 2021.
3. Lu, Tianjian, Thibault Marin, Yue Zhuo, Yi-Fan Chen, and Chao Ma. "Accelerating MRI Reconstruction on TPUs." In *2020 IEEE High Performance Extreme Computing Conference (HPEC)*, pp. 1-9. IEEE, 2020.

- DFT is critical in many scientific and engineering applications.
- General form of DFT

$$X_k \triangleq X(z_k) = \sum_{n=0}^{N-1} x_n z_k^{-n}$$

- Matrix Form

$$\{X\} = [V] \{x\},$$

where

$$\begin{aligned} \{X\} &= (X(z_0), X(z_1), \dots, X(z_{N-1}))^T, \\ \{x\} &= (x_0, x_1, \dots, x_{N-1})^T, \end{aligned}$$

and

$$[V] = \begin{pmatrix} 1 & z_0^{-1} & z_0^{-2} & \cdots & z_0^{-(N-1)} \\ 1 & z_1^{-1} & z_1^{-2} & \cdots & z_1^{-(N-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & z_{N-1}^{-1} & z_{N-1}^{-2} & \cdots & z_{N-1}^{-(N-1)} \end{pmatrix}.$$

- Three-dimensional (3D) DFT

$$X(z_{1k}, z_{2k}, z_{3k}) = \sum_{n_1=0}^{N_1-1} \sum_{n_2=0}^{N_2-1} \sum_{n_3=0}^{N_3-1} x(n_1, n_2, n_3) z_{1k}^{-n_1} z_{2k}^{-n_2} z_{3k}^{-n_3}$$

- Matrix Form

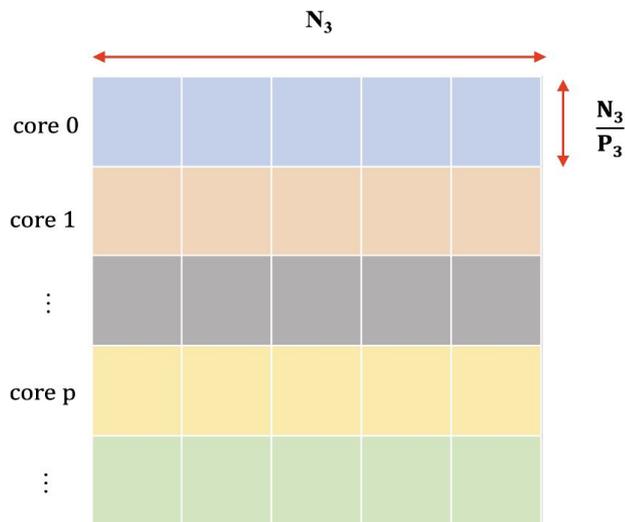
$$\{X\} = [V_1] \otimes [V_2] \otimes [V_3] \{x\},$$

where

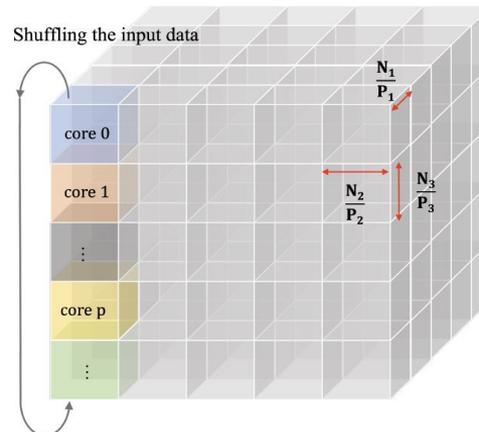
$$[V_j] = \begin{pmatrix} 1 & z_{j0}^{-1} & z_{j0}^{-2} & \cdots & z_{j0}^{-(N_j-1)} \\ 1 & z_{j1}^{-1} & z_{j1}^{-2} & \cdots & z_{j1}^{-(N_j-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & z_{j,N_j-1}^{-1} & z_{j,N_j-1}^{-2} & \cdots & z_{j,N_j-1}^{-(N_j-1)} \end{pmatrix}$$

$j \in \{1, 2, 3\}.$

- Advantages of the one-shuffle scheme:
 - tensor operations are **localized** on individual cores;
 - communication (sending and receiving data among cores) takes place **along the same direction** on the interconnect network;
 - and it achieves **high parallel efficiency**.



(a)



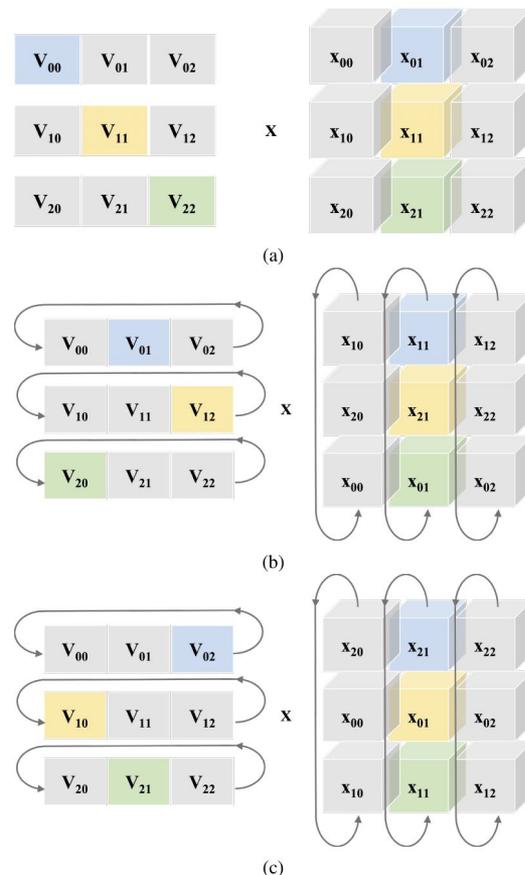
(b)

Algorithm 1 The one-shuffle scheme

```

1: function ONE_SHUFFLE( $v$ ,  $x$ ,  $\text{core\_idx}$ ,  $\text{num\_cores}$ ,
    $\text{src\_tgt\_pairs}$ )
2:    $\text{iteration\_idx} \leftarrow 0$ 
3:    $\text{slice\_idx} \leftarrow \text{core\_idx}$ 
4:    $x_{\text{out}} \leftarrow \text{einsum}(v[\text{slice\_idx}], x)$ 
5:    $\text{slice\_idx} \leftarrow \text{mod}(\text{slice\_idx} + 1, \text{num\_cores})$ 
6:   while  $\text{iteration\_idx} < \text{num\_cores} - 1$  do
7:      $x \leftarrow \text{collective\_permute}(x, \text{src\_tgt\_pairs})$ 
8:      $x_{\text{out}} \leftarrow x_{\text{out}} + \text{einsum}(v[\text{slice\_idx}], x)$ 
9:      $\text{slice\_idx} \leftarrow \text{mod}(\text{slice\_idx} + 1, \text{num\_cores})$ 
10:     $\text{iteration\_idx} \leftarrow \text{iteration\_idx} + 1$ 
11:   return  $x_{\text{out}}$ 

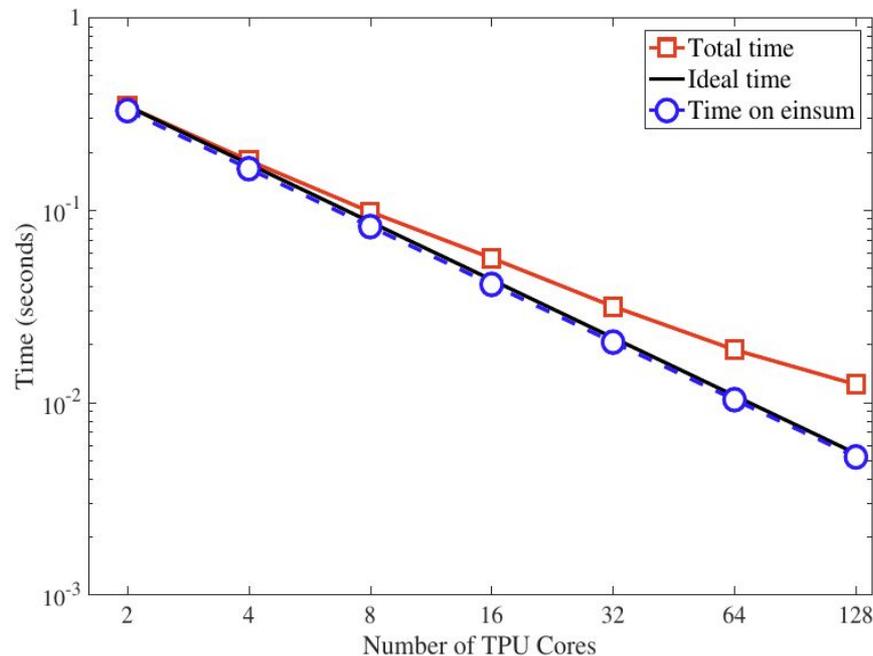
```



- Strong scaling analysis
- Define ideal time of linear scaling as a reference

$$\text{ideal time} = \frac{T_2}{\frac{N_{\text{core}}}{2}}$$

- The total computation time has a **close-to-linear scaling**.

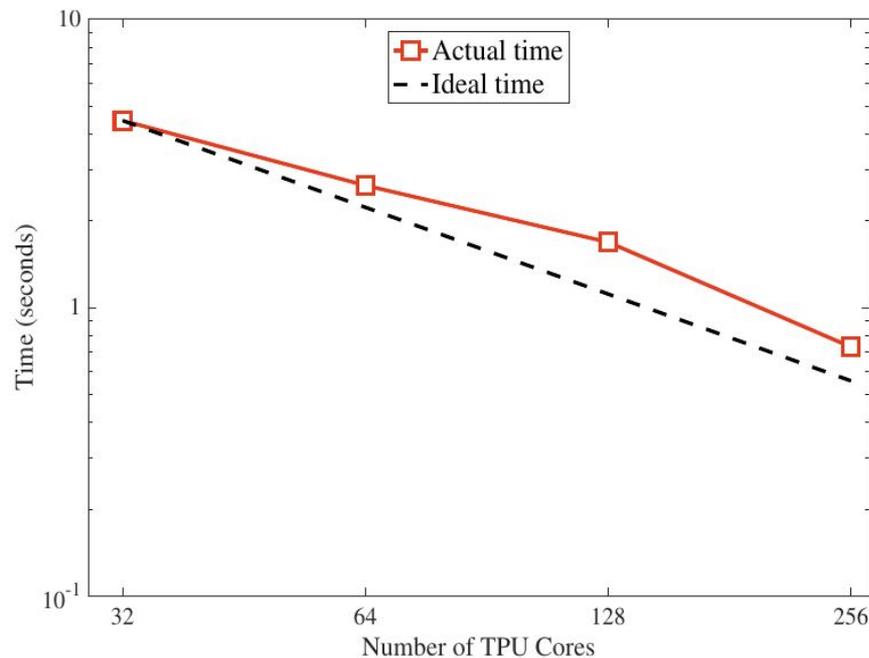


2D DFT, fixed problem size of 8192 x 8192;
up to 128 TPU cores being used.

- Strong scaling analysis:
- Define ideal time of linear scaling as a reference

$$\text{ideal time} = \frac{T_{32}}{\frac{N_{\text{core}}}{32}}$$

- The total computation time has a **close-to-linear scaling**.



3D DFT, fixed problem size of 2048 x 2048 x 2048, and up to 256 TPU cores being used.

- The FFT formulation starts with

$$X_k \triangleq \sum_{n=0}^{N-1} x_n e^{-j2\pi \frac{nk}{N}}$$

- The global index n can be expressed as

$$n = Pl + \beta,$$

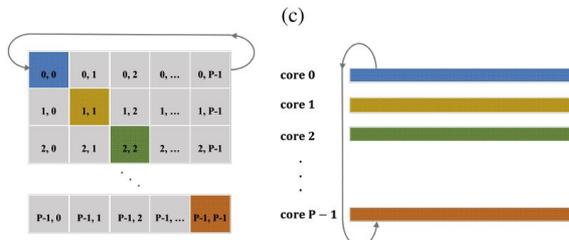
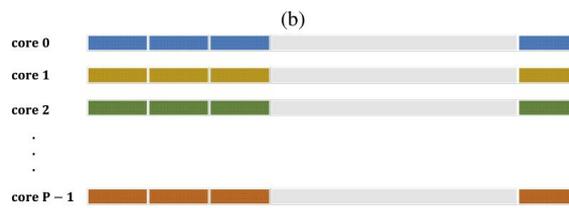
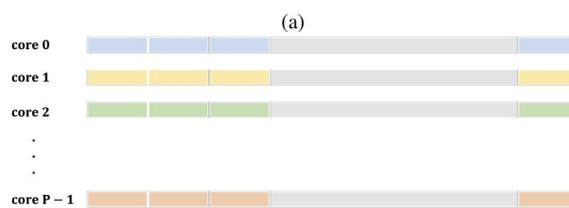
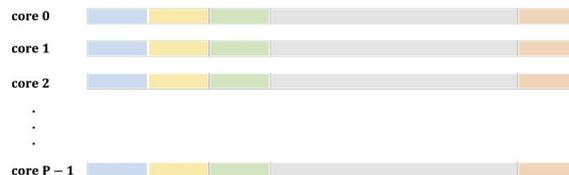
where $l = 0, 1, \dots, \frac{N}{P} - 1$ and $\beta = 0, 1, \dots, P - 1$.

- Rewrite as **phase adjustment** and **localized transform**

$$\begin{aligned} X_k &\triangleq \sum_{n=0}^{N-1} x_{(Pl+\beta)} e^{-j2\pi \frac{(Pl+\beta)k}{N}} \\ &= \sum_{\beta=0}^{P-1} e^{-j2\pi \frac{\beta k}{N}} \left(\sum_{l=0}^{\frac{N}{P}-1} x_{(Pl+\beta)} e^{-j2\pi \frac{lk}{\frac{N}{P}}} \right). \end{aligned}$$

The parallel algorithm

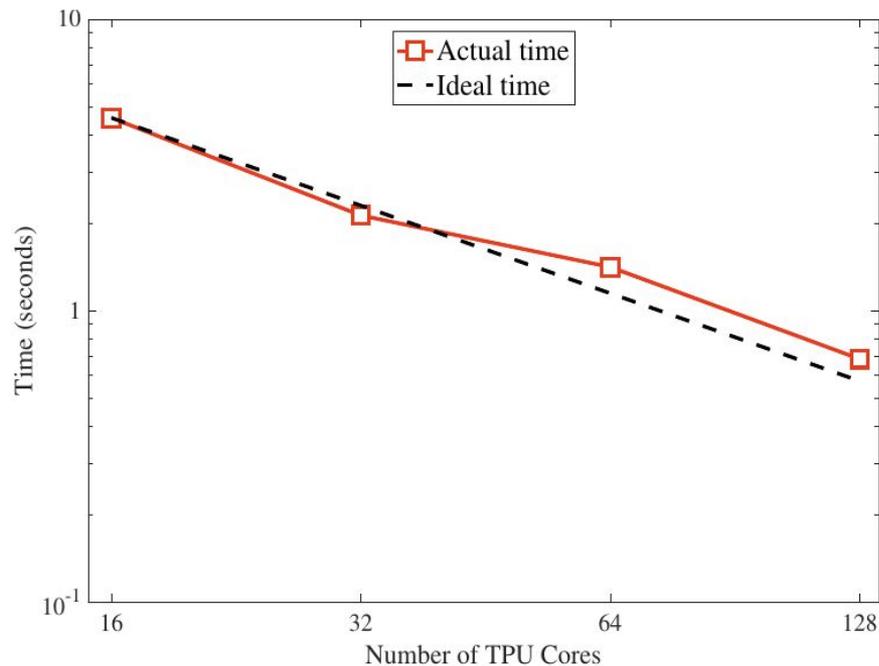
- (a) Data decomposition.
- (b) The gathering of input for in-order transform.
- (c) The transform performed locally on individual cores.
- (d) Applying phase adjustment with the **one-shuffle algorithm**.



- Strong scaling analysis.
- Define ideal time of linear scaling as a reference

$$\text{ideal time} = \frac{T_{16}}{\frac{N_{\text{core}}}{16}}$$

- The total computation time has a **close-to-linear scaling**.

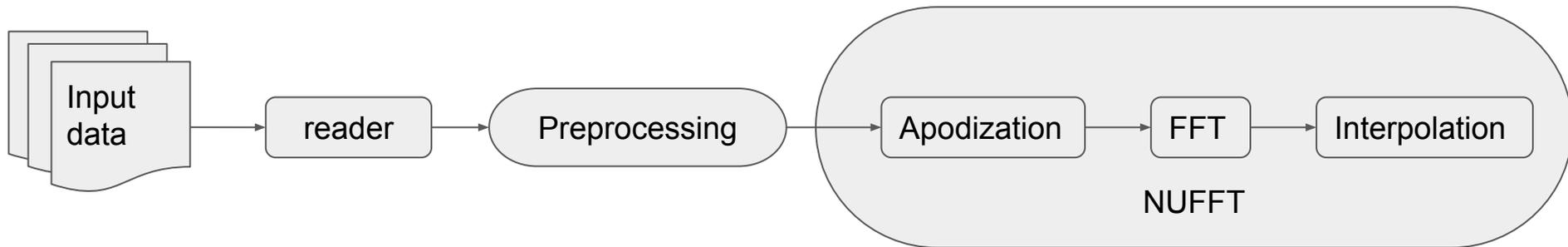


3D FFT, fixed problem size of 2048 x 2048 x 2048, and up to 128 TPU cores being used.

- The computation time of a few 3D DFT and FFT examples on a full TPU v3 Pod with 2048 cores.
- The runtimes reported in the table are for **complex** transforms.
- As a reference, the runtime of a **real** FFT for the problem size $8192 \times 8192 \times 8192$ on **CPUs**: 2048 nodes of Fujitsu PRIMERGY CX1640 M1 cluster is 5.36 seconds (converted from 10 TFlops, D. Takahashi 2019).

No.	Problem size	Time (seconds)	
		DFT	FFT
1	$8192 \times 8192 \times 8192$	12.66	8.30
2	$4096 \times 4096 \times 4096$	1.07	1.01
3	$2048 \times 2048 \times 2048$	0.120	0.118
4	$1024 \times 1024 \times 1024$	0.0220	0.0148

D. Takahashi, "Implementation of parallel 3-D real FFT with 2-D decomposition on Intel Xeon Phi clusters," in *International Conference on Parallel Processing and Applied Mathematics*. Springer, 2019, pp. 151–161.



- NUFFT

$$s(k_{x,m}, k_{y,m}) = \sum_{n=1}^N \rho_n e^{-i2\pi(k_{x,m}x_n + k_{y,m}y_n)}$$

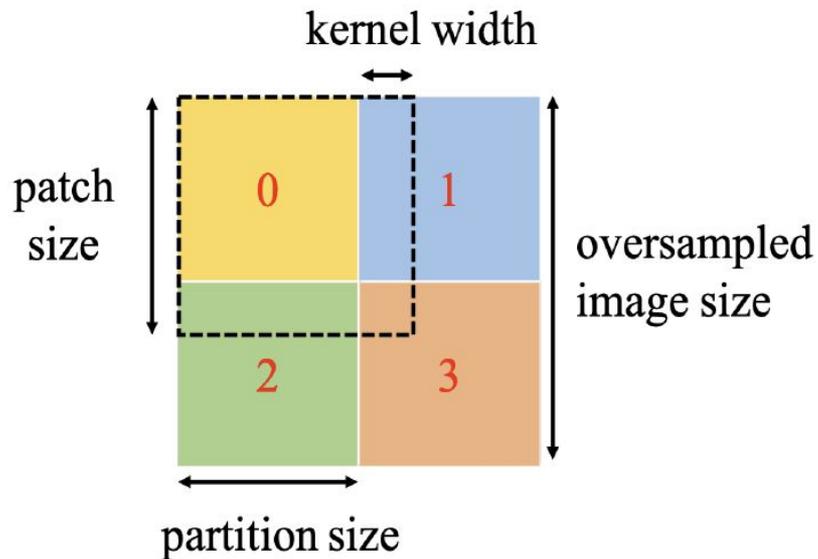
where $(k_{x,m}, k_{y,m})$, $m = 1, 2, \dots, M$ represents the k -space coordinates on a nonuniform grid, (x_n, y_n) , $n = 1, 2, \dots, N$ represents the spatial coordinates on a uniform grid, and ρ_n denotes the image intensity on grid (x_n, y_n) .

- Matrix form

$$\mathbf{s} = \mathbf{C}\mathbf{F}\mathbf{D}\boldsymbol{\rho}$$

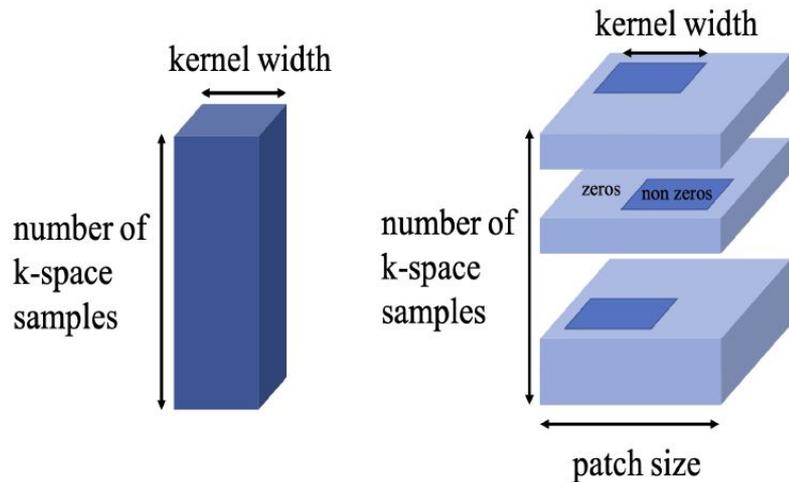
where \mathbf{D} is the apodization operator, \mathbf{F} denotes the FFT operator, and \mathbf{C} represents the interpolation operator.

Preprocessing

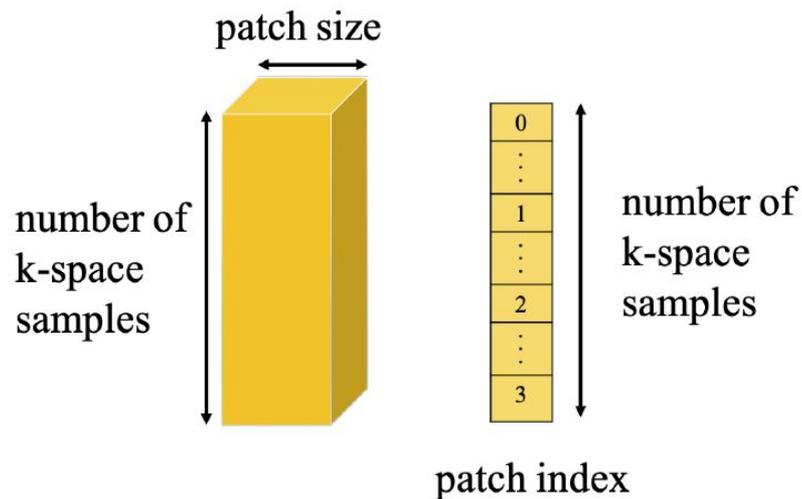


Perform checker-board partition to an oversampled image into patches.

Preprocessing

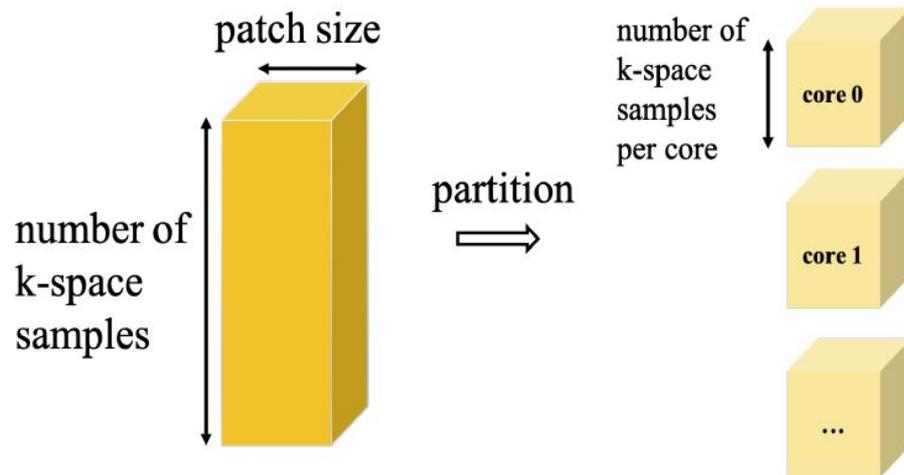


Pad tensors of kernel coefficients with zeros.



Shuffle kernel coefficients along the k-space dimension.

Load balancing

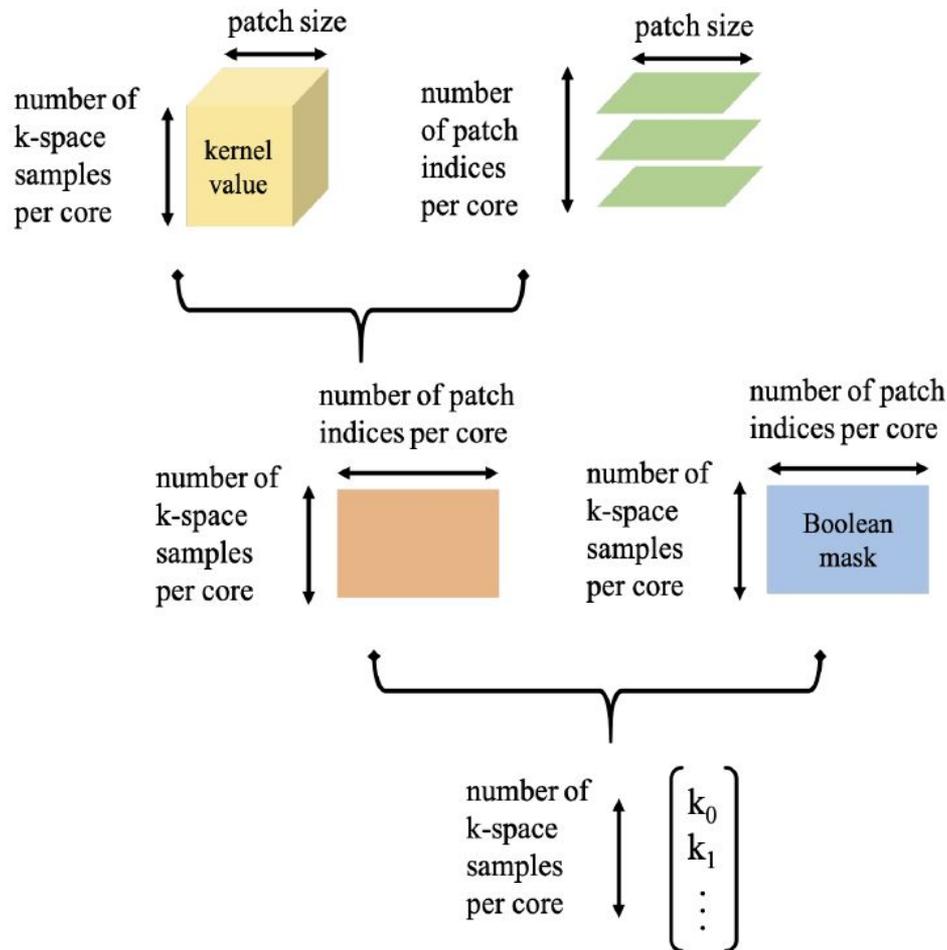


Each TPU core contains partial k-space information.

Interpolation in Forward NUFFT

Transform onto nonuniform grids:

- **Tensor contraction** between kernel coefficients and patches along patch dimension.
- **Matrix multiplication** with a Boolean mask.
- **Reduce sum** along the patch dimension.



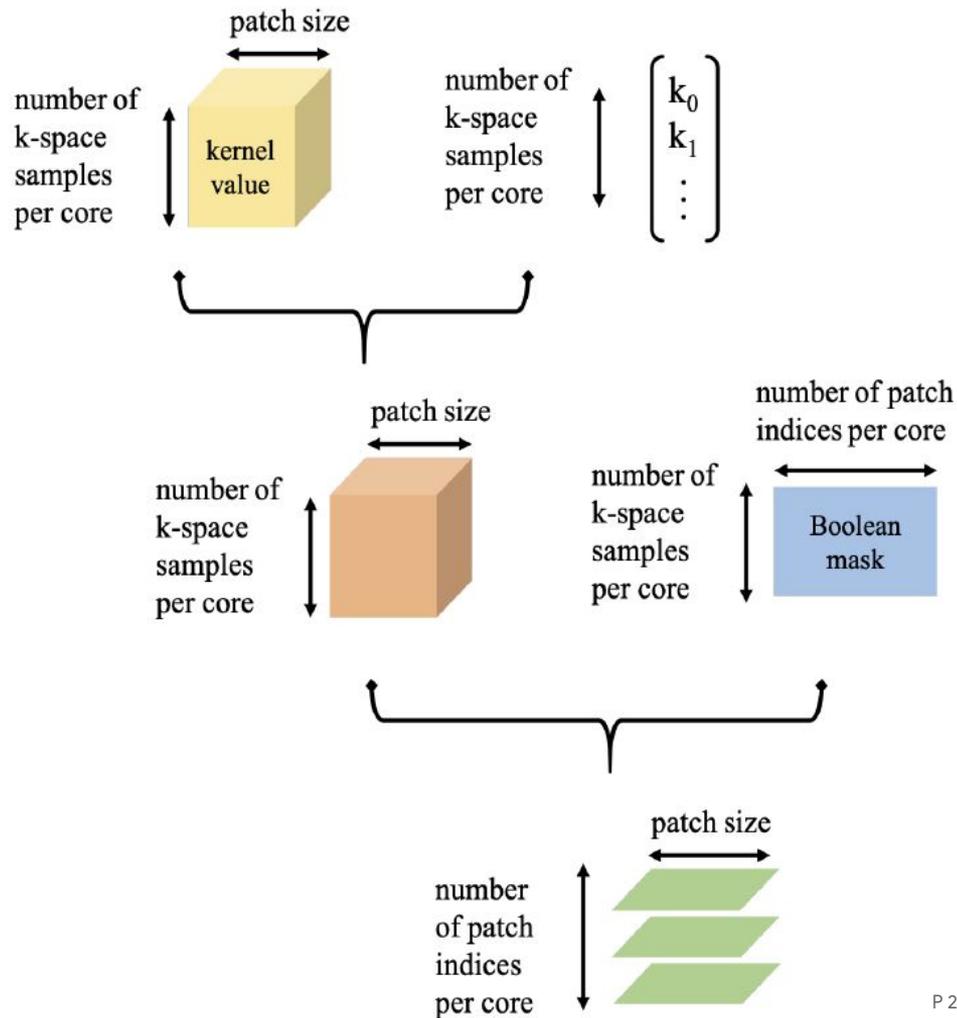
Case Study 3: NUFFT on TPUs

Interpolation in Adjoint NUFFT

Transform onto the uniform grids:

- **Scale** the kernel coefficients with k-space data
- **Tensor contraction** with a Boolean mask along the k-space dimension.

Patches are assembled back to the image.



Computation on three types of hardware for forward NUFFT:

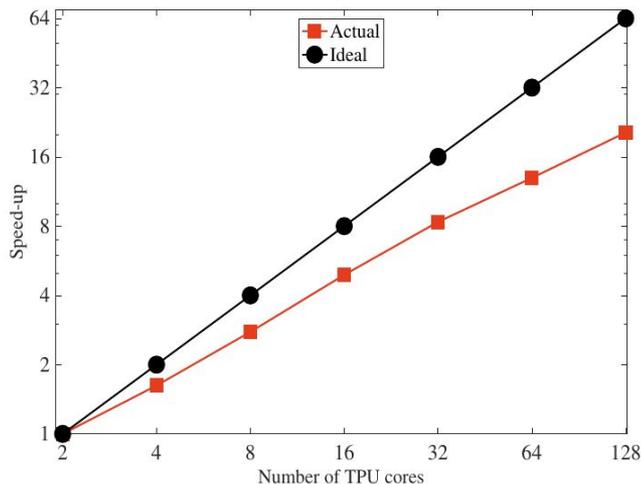
- CPU: Intel(R) Xeon(R) Silver 4110
8-core 2.10 GHz
- GPU: Nvidia V100
- TPU: one TPU v3 unit (eight cores).

Image size	Time (ms)		
	CPU	GPU	TPU
64 x 64	18.97	2.89	0.11
128 x 128	75.04	3.08	0.36
256 x 256	250.5	3.02	1.23
512 x 512	1135.19	3.04	5.41

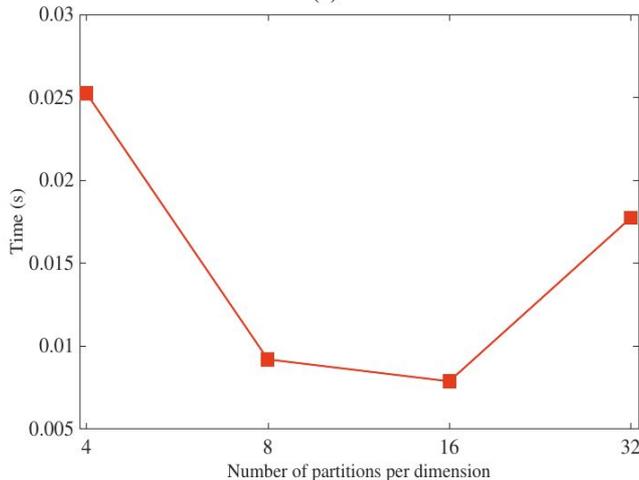
Case Study 3: NUFFT on TPUs

Adjoint NUFFT

- Image size: 512 x 512
- Oversampling factor: 2
- Number of points in k-space: 412,160
- Strong scaling:
 - Number of TPU cores: 2 to 128
 - Number of partitions: 16 along each dimension
- Computation time versus partitions
 - 16 TPU cores



(a)



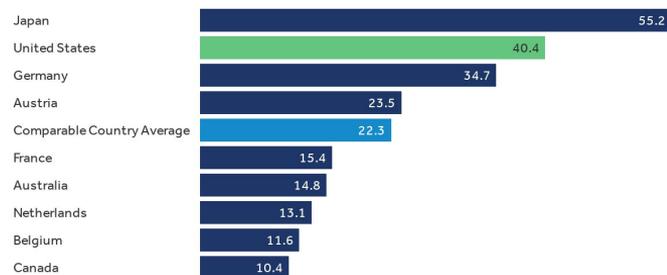
(b)

Magnetic resonance imaging (MRI) is a **powerful** medical imaging modality:

- non-invasive
- excellent soft-tissue contrast
- high spatial resolution

MRI has revolutionized the field of medical imaging since its invention in 1970s.

Magnetic Resonance Imaging (MRI) units per million population, 2019



Notes: Data for Austria and the Netherlands are from 2018. Data for Germany and Japan are from 2017.

Source: Kaiser Family Foundation Analysis of OECD Data

Peterson, KFF
Health System Tracker

Computation in MR image reconstruction is now the new **bottleneck**:

- MR data acquisition speed is approaching the physical limits.
- Further acceleration of MR requires breaking the Nyquist sampling criterion by sparse sampling and constrained image reconstruction.
- However, the state-of-the-art MR image reconstruction methods often build upon **large-scale, iterative, optimization** algorithms
 - with extensive usage of **non-uniform Fourier transform**
 - **computationally infeasible** for practical clinical use.

- MRI signal model

$$d_{\kappa,\gamma} = \sum_n s_{n,\gamma} \rho_n e^{-i2\pi \mathbf{k}_\kappa \cdot \mathbf{r}_n}$$

$$= [\mathbf{F}(\boldsymbol{\rho})]_{\kappa,\gamma}$$

- In compressed sensing, one reconstructs an image from the undersampled k-space data by solving

$$\min_{\boldsymbol{\rho}} \|\mathbf{F}(\boldsymbol{\rho}) - \mathbf{d}\|_2^2 + \lambda \|\boldsymbol{\Theta}(\boldsymbol{\rho})\|_1$$

Data fidelity

Sparsity
constraint

- We use the Alternating Direction Method of Multipliers (**ADMM**) to solve the large-scale convex optimization problem

$$\min_{\boldsymbol{\rho}} \|\mathbf{F}(\boldsymbol{\rho}) - \mathbf{d}\|_2^2 + \lambda \|\boldsymbol{\mu}\|_1$$

$$\text{s.t. } \boldsymbol{\Theta}(\boldsymbol{\rho}) - \boldsymbol{\mu} = 0$$

- ADMM consists of three updates:

$$\boldsymbol{\mu}^{m+1} = \operatorname{argmin}_{\boldsymbol{\mu}} \lambda \|\boldsymbol{\mu}\|_1 + \frac{\beta}{2} \|\boldsymbol{\Theta}(\boldsymbol{\rho}^m) - \boldsymbol{\mu} + \boldsymbol{\eta}^m\|_2^2$$

$$\boldsymbol{\rho}^{m+1} = \operatorname{argmin}_{\boldsymbol{\rho}} \|\mathbf{F}(\boldsymbol{\rho}) - \mathbf{d}\|_2^2 + \frac{\beta}{2} \|\boldsymbol{\Theta}(\boldsymbol{\rho}) - \boldsymbol{\mu}^{m+1} + \boldsymbol{\eta}^m\|_2^2$$

$$\boldsymbol{\eta}^{m+1} = \boldsymbol{\eta}^m + \boldsymbol{\Theta}(\boldsymbol{\rho}^{m+1}) - \boldsymbol{\mu}^{m+1}$$

- The update of the auxiliary variable has a closed-form solution

$$\boldsymbol{\mu}^{m+1} = S_{\frac{\lambda}{\beta}} (\boldsymbol{\Theta} (\boldsymbol{\rho}^m) + \boldsymbol{\eta}^m)$$

and the element-wise soft thresholding can be written as

$$S_{\frac{\lambda}{\beta}}(\alpha) = \begin{cases} \alpha - \frac{\lambda}{\beta}, & \alpha > \frac{\lambda}{\beta} \\ 0, & |\alpha| \leq \frac{\lambda}{\beta} \\ \alpha + \frac{\lambda}{\beta}, & \alpha < -\frac{\lambda}{\beta} \end{cases}$$

- The update of the primal variable (complex image intensities) can be considered as a regularized least square problem.
- The necessary and sufficient optimality condition is

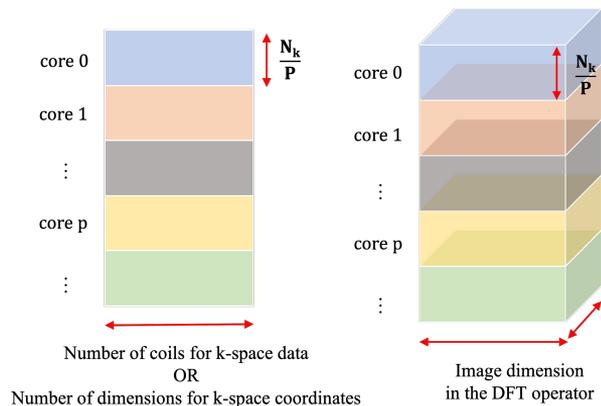
$$\mathbf{A}\boldsymbol{\rho}^{m+1} = \mathbf{b}$$

where

$$\mathbf{A} = \mathbf{F}^H \mathbf{F} + \frac{\beta}{2} \boldsymbol{\Theta}^H \boldsymbol{\Theta},$$

$$\mathbf{b} = \mathbf{F}^H \mathbf{d} + \frac{\beta}{2} \boldsymbol{\Theta}^H (\boldsymbol{\mu}^{m+1} - \boldsymbol{\eta}^m)$$

- This is solved iteratively by using the **conjugate gradient (CG)** method.



Data decomposition applied to data and DFT operator.

- The data decomposition is applied to the k-space.
- DFT and sparsifying transform operations
 - The DFT operation and its adjoint are formulated as **tensor contractions** (`tf.einsum`).
 - The sparsifying transform operation and its adjoint are formulated as **convolutions** (`tf.nn.conv1d`).

Algorithm 1 The generation of the DFT operator on TPUs

```

1: function MAP_TO_UNIT_CIRCLE(k_r_product)
2:   return  $\exp(-i 2\pi k_r\_product)$ 

1: function GEN_DFT_OPERATOR(k_coord, image_coord)
2:   kr_dim0  $\leftarrow$  map_to_unit_circle(
       einsum(image_coord[0], k_coord[:,0]))
3:   kr_dim1  $\leftarrow$  map_to_unit_circle(
       einsum(image_coord[1], k_coord[:,1]))
4:   dft_op  $\leftarrow$  einsum(kr_dim0, kr_dim1)
5:   return dft_op
    
```

- Communication strategy
- ADMM has three updates per iteration:
 - The update of the **auxiliary** variable is **local**.
 - The update of the **dual** variable is **local**.
 - The update of the **primal** variable is through CG solver, **requiring communication** (tf.cross_replica_sum) to sum the partial images across TPU cores such that all cores start the new CG iteration with the same image.

Algorithm 4 ADMM on TPUs

```

1: function ADMM_STEP( $\eta, \rho, \rho_0$ )
2:    $\mu_{next} \leftarrow \text{update\_auxiliary\_var}(\rho, \eta)$ 
3:    $\rho_{next} \leftarrow \text{update\_primal\_var}(\mu_{next}, \rho_0, \eta)$ 
4:    $\alpha \leftarrow \text{get\_relative\_diff}(\rho_{next}, \rho)$ 
5:    $\eta_{next} \leftarrow \text{update\_dual\_var}(\mu_{next}, \rho, \eta)$ 
6:   return  $\eta_{next}, \rho_{next}, \alpha$ 

1: function ADMM_RECONSTRUCT( $\rho_0, \text{max\_iterations}, \text{rtol}$ )
2:    $\triangleright \rho_0$  contains initial values of  $\rho$  and  $\text{rtol}$  is the relative
   tolerance in terms of the squared norm of the residual.
3:    $\eta, \rho \leftarrow \text{get\_initial\_value}(\rho_0)$ 
4:    $i \leftarrow 0$ 
5:    $\alpha \leftarrow 1.0$ 
6:   while  $i < \text{max\_iterations} \ \& \ \alpha > \text{square}(\text{rtol})$  do
7:      $\eta, \rho, \alpha \leftarrow \text{admm\_step}(\eta, \rho, \rho_0)$ 
8:      $i \leftarrow i + 1$ 
9:   return  $\rho$ 
    
```

Algorithm 2 The conjugate gradient method

```

1: function CG_STEP( $\text{linear\_op}, r, d, x, \tau$ )  $\triangleright d$  is the
   conjugate vector and  $r$  is the the residual vector.
2:    $a\_d \leftarrow \text{linear\_op}(d)$ 
3:    $\alpha \leftarrow \text{divide}(\tau, \text{dot\_product}(d, a\_d))$ 
4:    $x_{next} \leftarrow x + \alpha d$ 
5:    $r_{next} \leftarrow r - \alpha a\_d$ 
6:    $\tau_{next} \leftarrow \text{dot\_product}(r_{next}, r_{next})$ 
7:    $\beta \leftarrow \text{divide}(\tau_{next}, \tau)$ 
8:    $d_{next} \leftarrow r_{next} + \beta d$ 
9:   return  $r_{next}, d_{next}, x_{next}, \tau_{next}$ 

1: function CONJUGATE_GRADIENT( $\text{linear\_op}, b, x_0$ ,
    $\text{max\_iterations}, \text{atol}$ )
2:    $\triangleright x_0$  contains initial values of  $x$  and  $\text{atol}$  is the absolute
   tolerance in terms of the norm of the residual vector.
3:    $x \leftarrow x_0$ 
4:    $r \leftarrow \text{linear\_op}(x)$ 
5:    $d \leftarrow r$ 
6:    $\tau \leftarrow \text{dot\_product}(r, r)$ 
7:    $i \leftarrow 0$ 
8:   while  $i < \text{max\_iterations} \ \& \ \tau > \text{square}(\text{atol})$  do
9:      $r_{next}, d_{next}, x_{next}, \rho_{next} \leftarrow \text{cg\_step}(\text{linear\_op}, r, d, x, \tau)$ 
10:     $i \leftarrow i + 1$ 
11:  return  $x$ 
    
```

Case Study 4: ADMM and CG on TPUs

Accuracy benchmark

- The k-space data were retrospectively undersampled with an undersampling factor of eight to demonstrate the capability of compressed sensing in accelerating MR.
- The total number of k-space measurements was 19,968 (1,664 samples per coil and 12 coils in total).
- The images were reconstructed on a 128 x 64 uniform grid.
- The relative difference is about **1%** for the voxels within the phantom, which is satisfactory.

ADMM				CG	
Regularization parameter	Augmented Lagrangian parameter	Relative tolerance	Maximum number of iterations	Absolute tolerance	Maximum number of iterations
1e-7	1.0	1e-4	5	1e-6	20

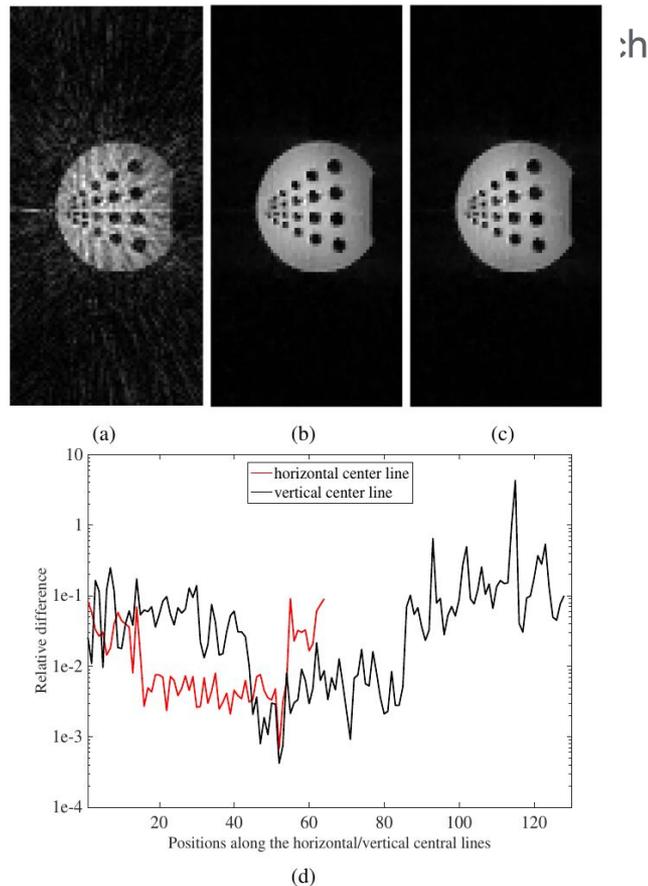


Fig. Reconstructed phantom images by using retrospectively undersampled data (a) with a single inverse DFT operation and (b) by the ADMM algorithm on CPU and (c) TPUs; and (d) the relative difference between the two images in (b) and (c) along the horizontal and vertical center lines of the image.

Case Study 4: ADMM and CG on TPUs

Parallel Efficiency

- The **strong scaling** analysis was adopted to understand the parallel efficiency.
- Phantom data were acquired by using 804 radial readouts, each with 1024 samples.
- The fully sampled k-space data were then retrospectively undersampled by a factor of eight, resulting in a total number of **1,241,076** k-space measurements (103,424 measurements per coil, 12 coils in total).
- Runtimes on CPU (**Intel(R) Xeon(R) Silver 4110** 8 core 2.1 GHz) and GPU (**NVIDIA V100** SXM2).

Hardware		Computation time (seconds)			
		CPU	GPU	TPU (number of TPU units)	
Non-uniform Fourier Transform		NUFFT	NUFFT	DFT	
Image Size	128 X 64	2.38	1.17	0.14 (1/4 unit)	0.036 (two units)
	1024 X 512	139.88	2.24	3.39 (16 units)	0.29 (128 units)

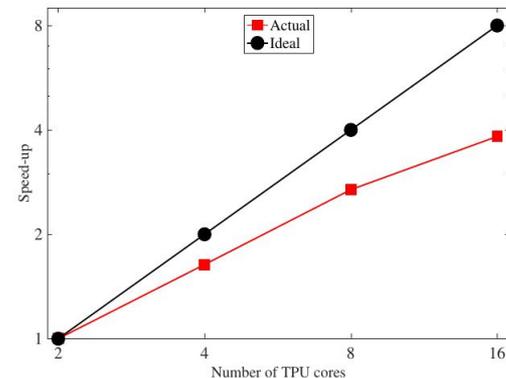


Fig. 5: The speed-up of reconstructing an image of size 128×64 with up to 16 TPU cores. The number of k -space measurements is 19,968 with 1,664 samples for each coil and 12 coils in total.

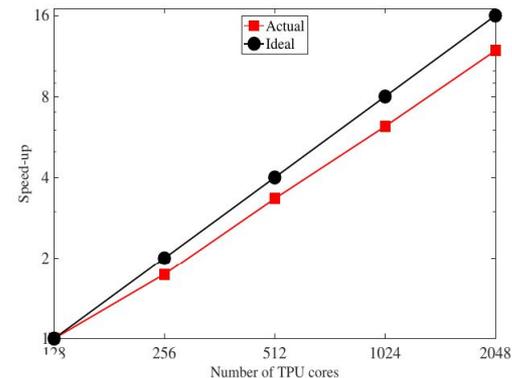


Fig. 6: The speed-up of reconstructing an image of size 1024×512 with up to 2048 TPU cores. The number of k -space measurements is 1,241,076 with 103,424 samples for each coil and 12 coils in total.

Conclusion

- Through the few case studies, we explore using TPU for scientific computing.
- The case studies include Fourier transform (DFT, FFT, NUFFT), linear system solver (CG), numerical optimization (ADMM), and their applications in medical imaging.
- We formulate the problem and design the algorithms in accordance with TPU's strength in tensor operations and its high-speed interconnect network.
- TPU achieves good acceleration for these scientific computing applications.

Thank You

tianjianlu@google.com

