# Deinsum: Practically I/O Optimal Multilinear Algebra

**Alexandros Nikolaos Ziogas**, **Grzegorz Kwasniewski**, Tal Ben-Nun, Timo Schneider, Torsten Hoefler
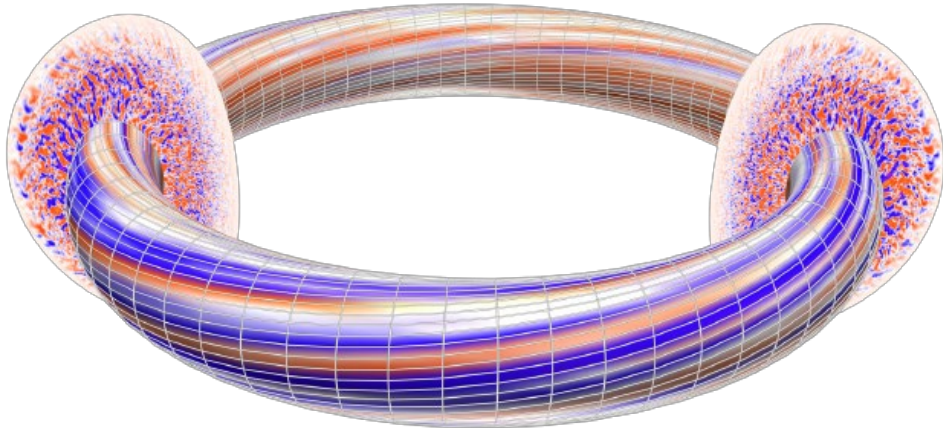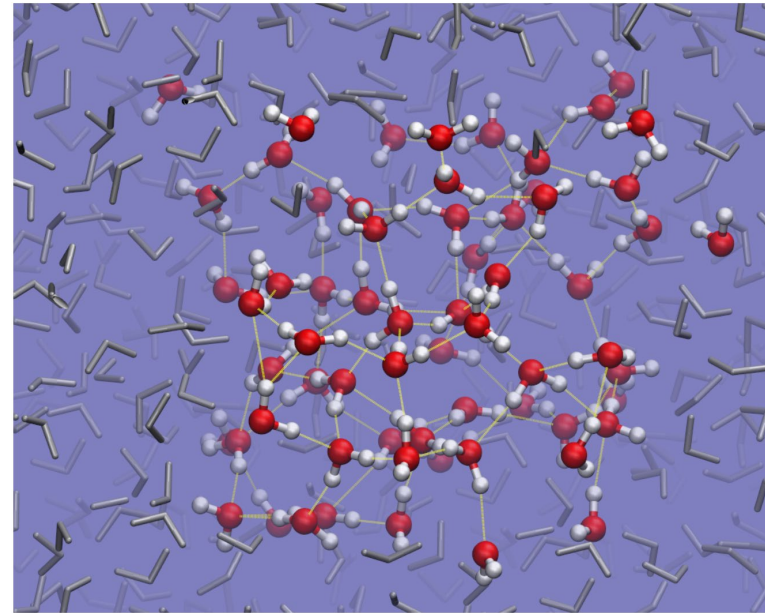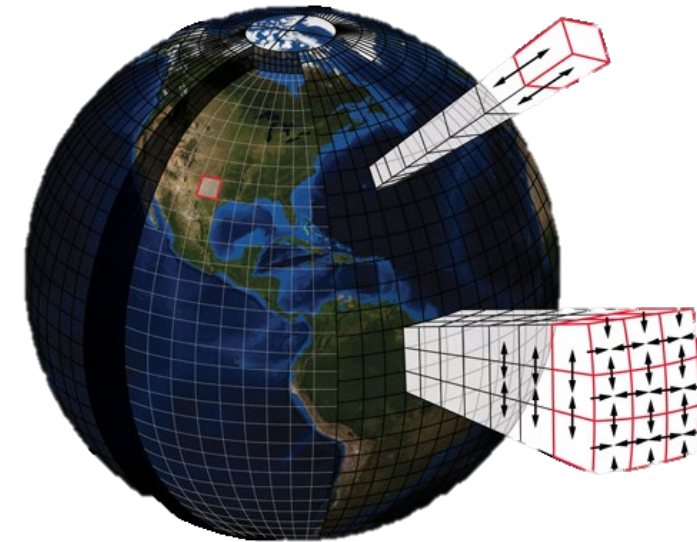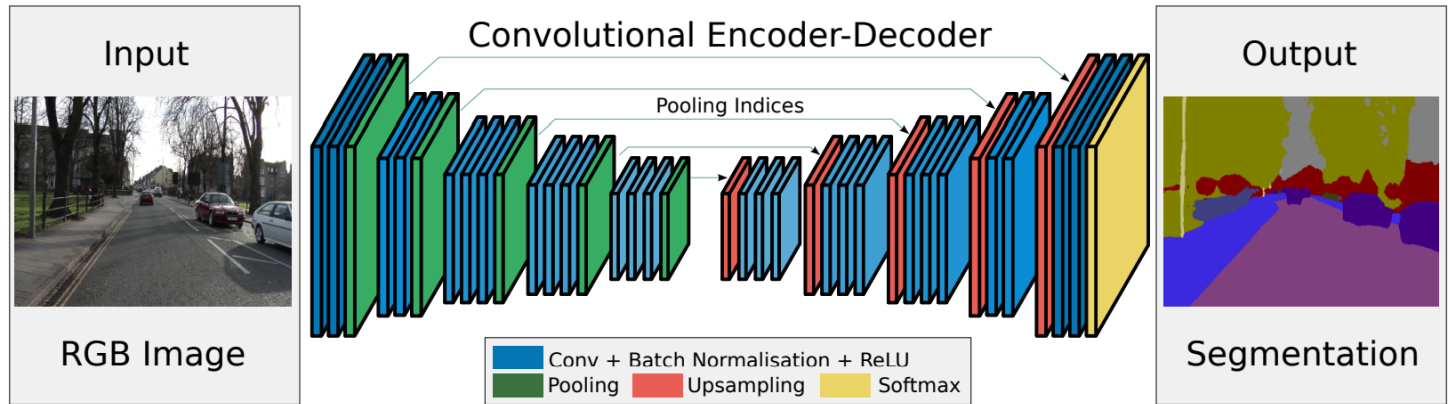
SC22

Dallas, TX | hpc accelerates.

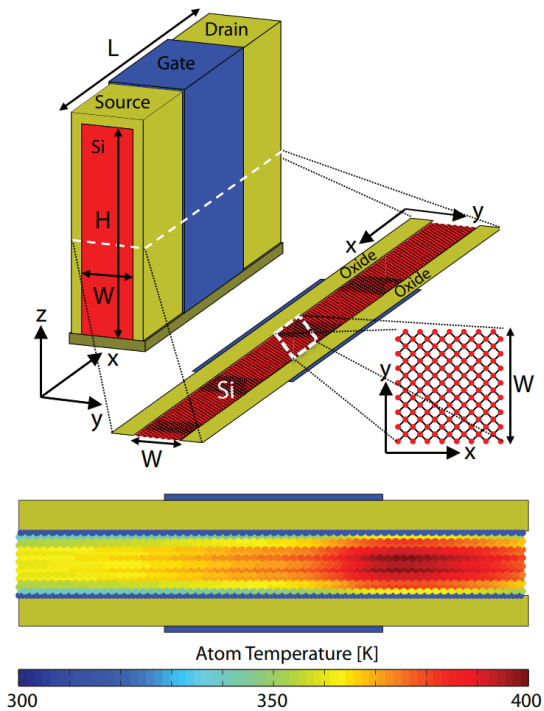Image generated from a GTS simulation by Kwan-Liu Ma and his group at the University of California, Davis.



Enabling Simulation at the Fifth Rung of DFT: Large Scale RPA Calculations with Excellent Time to Solution, Mauro Del Ben et al.



Credit: K. Cantner, AGI.



Atom Temperature [K]

300    350    400



SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation, Badrinarayanan et al.

Credit: Princeton Plasma Physics Laboratory

200 km

10 km

Credit: NASA's Goddard Space Flight Center

Credit: Jason Allen

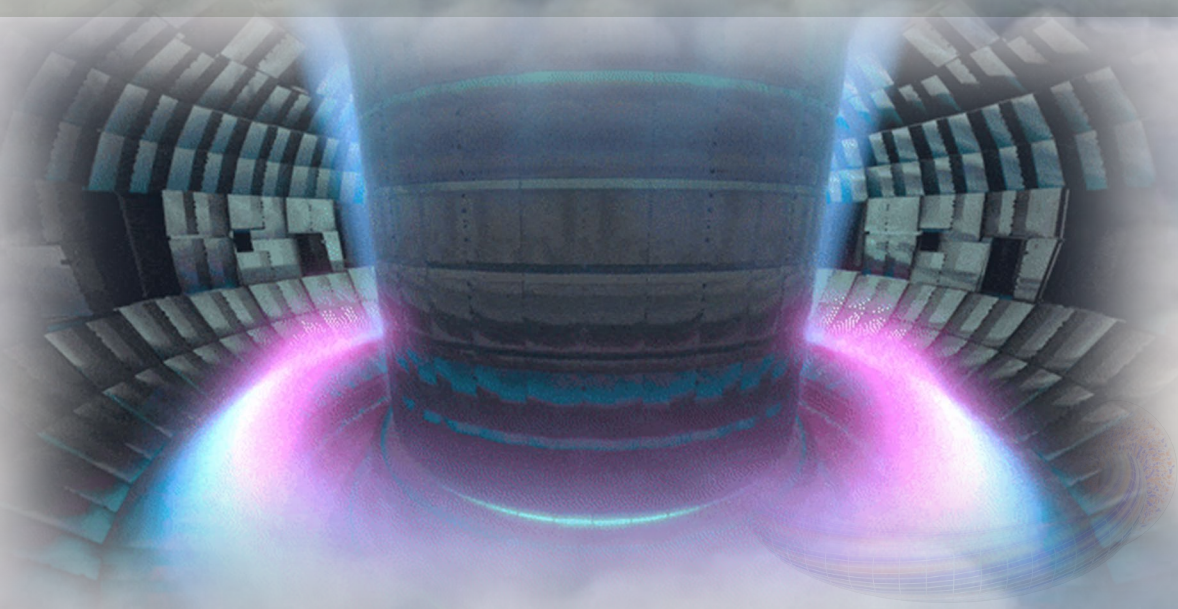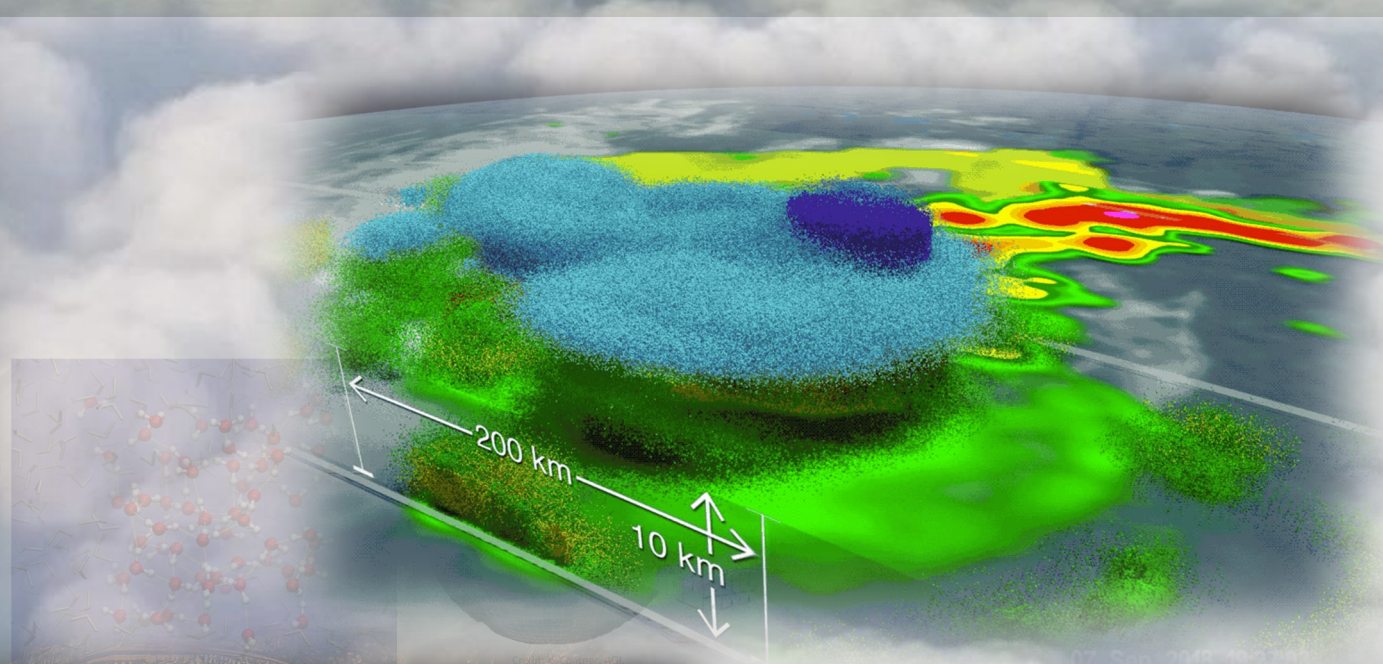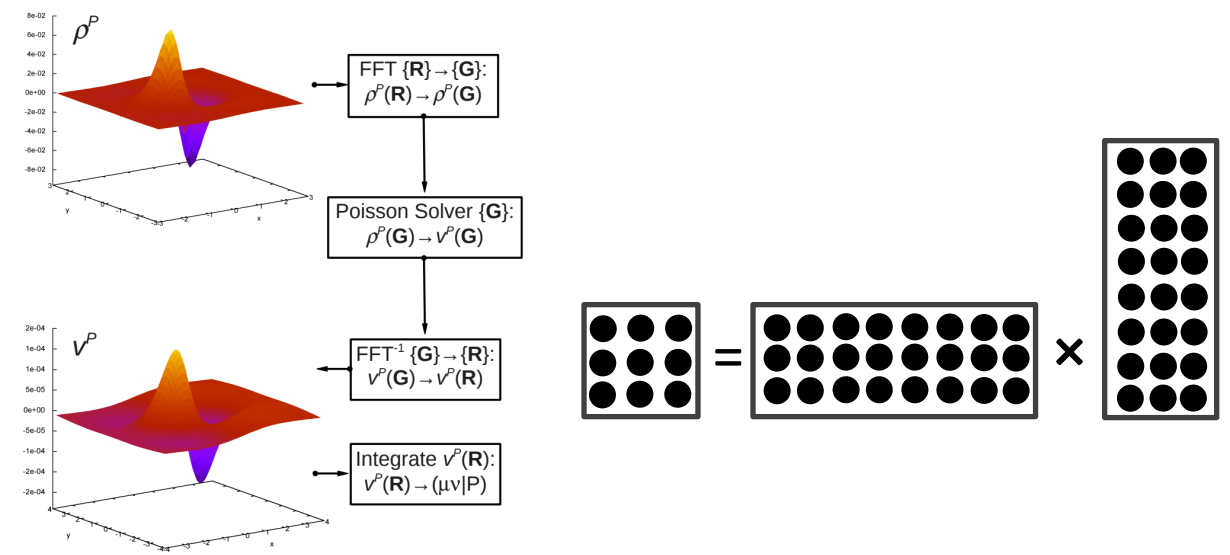Extreme Scale Plasma Turbulence Simulations on Top Supercomputers Worldwide, Tang et al.

Enabling Simulation at the Fifth Rung of DFT: Large Scale RPA Calculations with Excellent Time to Solution, Del Ben et al.

**NEGF**  SSE $\Sigma[G(E + \hbar\omega, k_z - q_z)\, D(\omega, q_z)](E, k_z)$

**Electrons** $G(E, k_z)$

$(E \cdot S - H - \Sigma^R) \cdot G^R = I$
$G^< = G^R \cdot \Sigma^< \cdot G^A$

**Phonons** $D(\omega, q_z)$

$(\omega^2 - \Phi - \Pi^R) \cdot D^R = I$
$D^< = D^R \cdot \Pi^< \cdot D^A$

GF

SSE  $\Pi[G(E, k_z)\, G(E + \hbar\omega, k_z + q_z)](\omega, q_z)$

A Data-Centric Approach to Extreme-Scale Ab initio Dissipative Quantum Transport Simulations, Ziogas et al.

$$ij, jk \rightarrow ik$$

*Einstein summation notation (einsum)*



```
for i in range(N):
    for j in range(N):
        for k in range(N):
            C[i,k]+=A[i,j]*B[j,k]
```

$$ij, jk \rightarrow ik$$

*Einstein summation notation (einsum)*

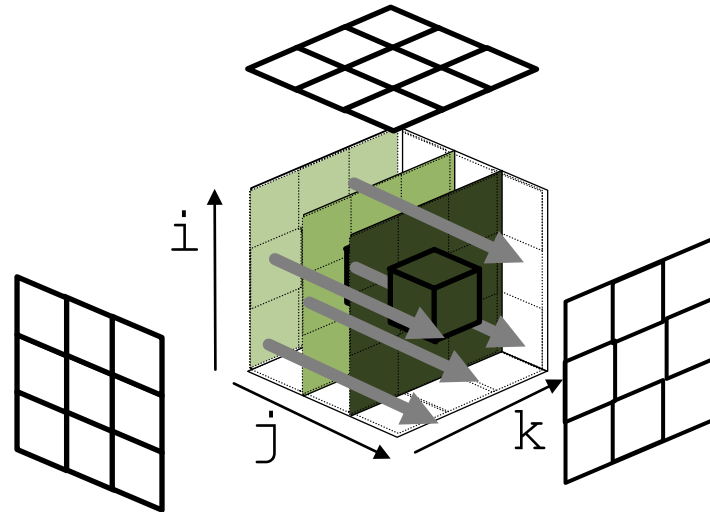| | |
|---|---|
| **Matrix-matrix multiplication** | $ij, jk \rightarrow ik$ |
| | $ij, jk, kl \rightarrow il$ |
| **Tensor and matrix chains** | $ij, jk, kl, lm \rightarrow im$ |
| | $ijk, ja, ka \rightarrow ia$ |
| | $ijk, ia, ka \rightarrow ja$ |
| | $ijk, ia, ja \rightarrow ka$ |
| | $ijklm, ja, ka, la, ma \rightarrow ia$ |
| **Long, higher order contraction chains** | $ijklm, ia, ja, la, ma \rightarrow ka$ |
| | $ijklm, ia, ja, ka, la \rightarrow ma$ |
| | $ijklm, jb, kc, ld, me \rightarrow ibcde$ |

Einsum

# Deinsum: Practically I/O Optimal Multi-Linear Algebra

**Simple Overlap Access Pattern (SOAP) data movement model**

**① Input**

$$ijk, ja, ka, al \rightarrow il$$

einsum, string

**② Split to binary operations**

$$ja, ka \rightarrow jka$$
$$ijk, jka \rightarrow ia$$
$$ia, al \rightarrow il$$

uses operation associativity to minimize #ops

**③ Communication-optimal schedules**

ja    ka    ijk    al

jka    ia    il

optimal tiling, parallel distribution, and kernel fusion

block-distributions

**④ Iteration spaces and distribution**

Iteration space partition:
`MPI_Cart_sub`

Distribute initial data:
`MPI_Broadcast`

i
j    k

**⑤ Automated code generation**

```
for j in range(NJ//P0J):
  for k in range(NK//P0K):
    for a in range(NA//P0A):
      t0[j, k, a] += A[j, a] * B[k, a]

t1 = np.tensordot(X, t0,
        axes=([1, 2], [0, 1]))
mpi.Allreduce(t1, comm=grid0_t1)

t2 = deinsum.Redistr
        comm1=grid0, comm=grid1)
```

auto generating DaCe-Python code compiled to shared library

**⑥ Results**

- **up to 19x speedup over SotA**
- **CPU and GPU support**

# Input

$$ijk, ja, ka, al \rightarrow il$$

```
for i in range(NI):
    for j in range(NJ):
        for k in range(NK):
            for l in range(NL):
                for a in range(NA):
                    out[i,l]+=X[i,j,k]*A[j,a]*B[k,a]*C[a,l]
```
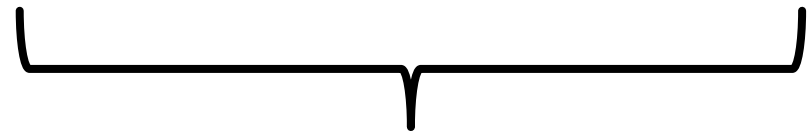
**Matrix-Matrix Product**

**naive implementation**
$4N_i N_j N_k N_l N_a$ **ops**

**Matricized Tensor Times Khatri-Rao Product (MTTKRP)**

**used in CANDECOMP/PARAFAC (CP) decomposition**

# Split to Binary Operations

②

$$ijk, ja, ka, al \rightarrow il$$

using Python module
`opt_einsum`

$$ja, ka \rightarrow jka$$
$$ijk, jka \rightarrow ia$$
$$ia, al \rightarrow il$$

# Split to Binary Operations

$$ja, ka \rightarrow jka$$

```
for j in range(NJ):
    for k in range(NK):
        for a in range(NA):
            t0[j,k,a]+=A[j,a]*B[k,a]
```

**Khatri-Rhao Product (KRP)**

$2N_j N_k N_a$ ops

$$ijk, jka \rightarrow ia$$

```
for i in range(NI):
    for j in range(NJ):
        for k in range(NK):
            for a in range(NA):
                t1[i,a]+=X[i,j,k]*t0[j,k,a]
```

```
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

$2N_i N_j N_k N_a$ ops

**Tensor DOT Product (TDOT)**

$$ia, al \rightarrow il$$

```
for i in range(NI):
    for l in range(NL):
        for a in range(NA):
            out[i,l]+=t1[i,a]*C[a,l]
```

```
out = t1 @ C
```

**Matrix-Matrix Product (GEMM)**

$2N_i N_l N_a$ ops

# Minimizing number of operations

②

```
for j in range(NJ):
    for k in range(NK):
        for a in range(NA):
            t0[j,k,a]+=A[j,a]*B[k,a]
```

```
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

```
out = t1 @ C
```

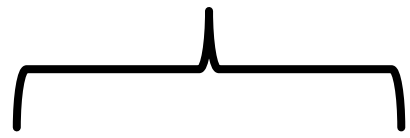```
for i in range(NI):
    for j in range(NJ):
        for k in range(NK):
            for l in range(NL):
                for a in range(NA):
                    out[i,l]+=(
                        X[i,j,k]*A[j,a]*
                        B[k,a]*C[a,l])
```

$2N_i N_j N_k N_a$
$+2N_j N_k N_a$
$+2N_i N_l N_a$ **ops**

**asymptotically**
$2N_i N_j N_k N_a$ **ops**

$<$

$4N_i N_j N_k N_l N_a$ **ops**

# Minimizing communication

③

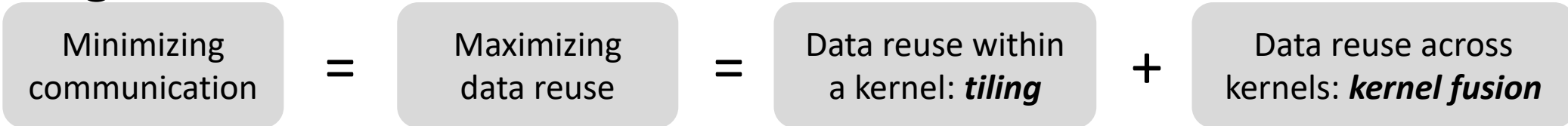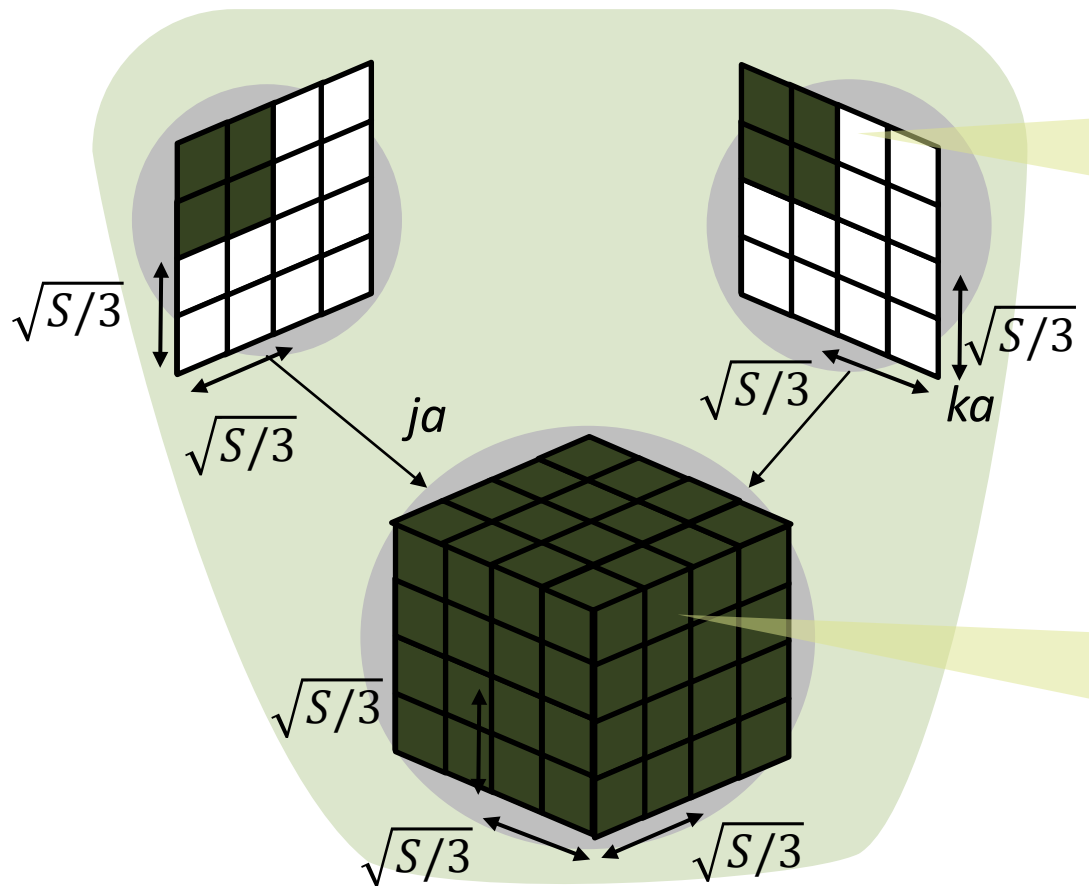| Minimizing communication | = | Maximizing data reuse | = | Data reuse within a kernel: *tiling* | + | Data reuse across kernels: *kernel fusion* |

$$ijk, ja, ka, al \rightarrow il \qquad ja, ka \rightarrow jka \qquad ijk, jka \rightarrow ia \qquad ia, al \rightarrow il$$

# Minimizing communication

③

| Minimizing communication | = | Maximizing data reuse | = | Data reuse within a kernel: *tiling* | + | Data reuse across kernels: *kernel fusion* |

$$ij k, ja, ka, al \rightarrow il \qquad ja, ka \rightarrow jka \qquad ij k, jka \rightarrow ia \qquad ia, al \rightarrow il$$



$\sqrt{S/3}$

$\sqrt{S/3}$  *ja*

$\sqrt{S/3}$  *ka*

$\sqrt{S/3}$

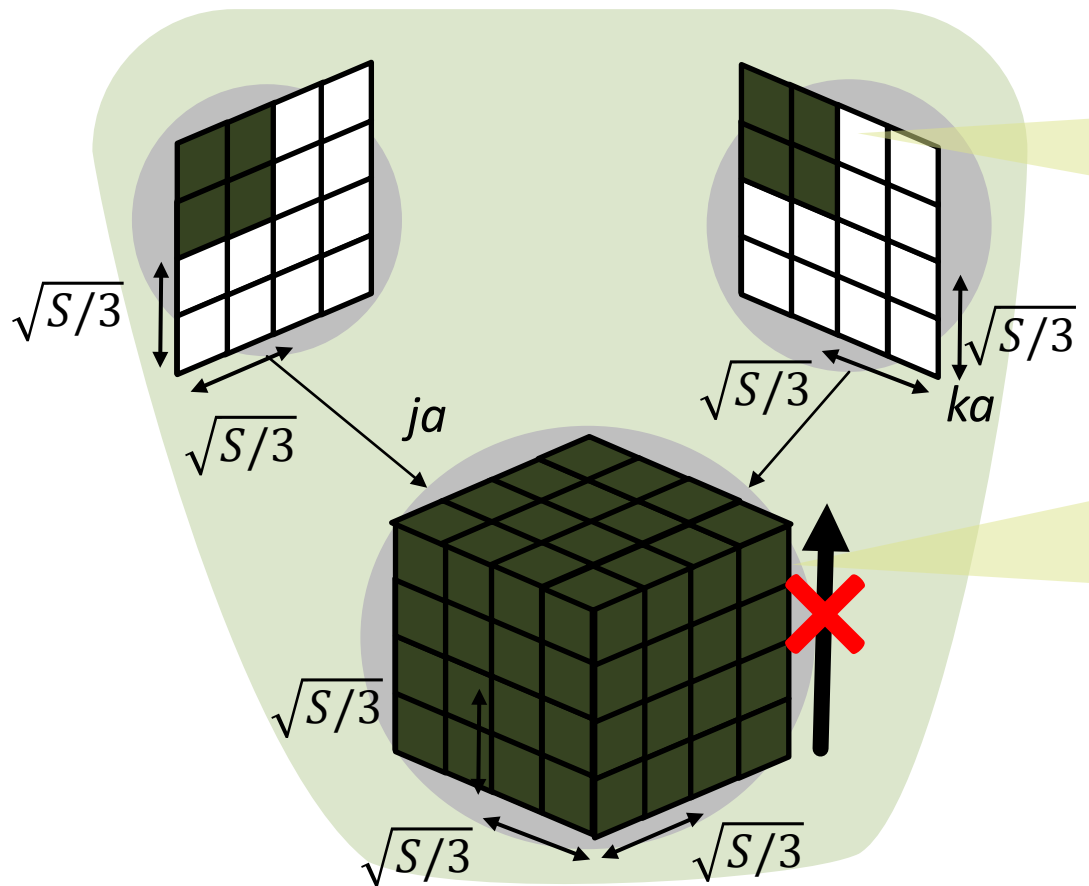$\sqrt{S/3}$

$\sqrt{S/3}$

$\sqrt{S/3}$

Inputs are too big to fit in local memory! Can store up to $S$ elements at once!

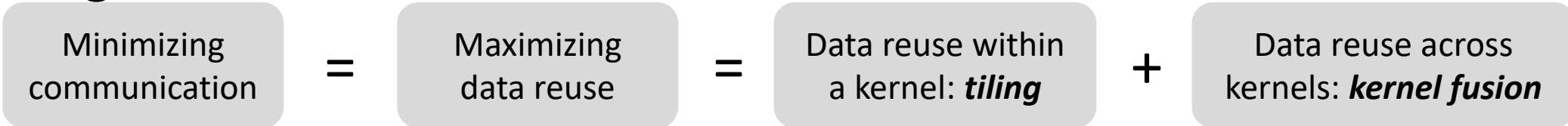*Data movement optimal classical matrix-matrix multiplication (SC'19)*

# Minimizing communication ③

| Minimizing communication | = | Maximizing data reuse | = | Data reuse within a kernel: *tiling* | + | Data reuse across kernels: *kernel fusion* |

$$ijk, ja, ka, al \rightarrow il \qquad \boxed{ja, ka \rightarrow jka} \quad ijk, jka \rightarrow ia \qquad ia, al \rightarrow il$$



$\sqrt{S/3}$

$\sqrt{S/3}$

*ja*

$\sqrt{S/3}$

*ka*

$\sqrt{S/3}$

$\sqrt{S/3}$

$\sqrt{S/3}$

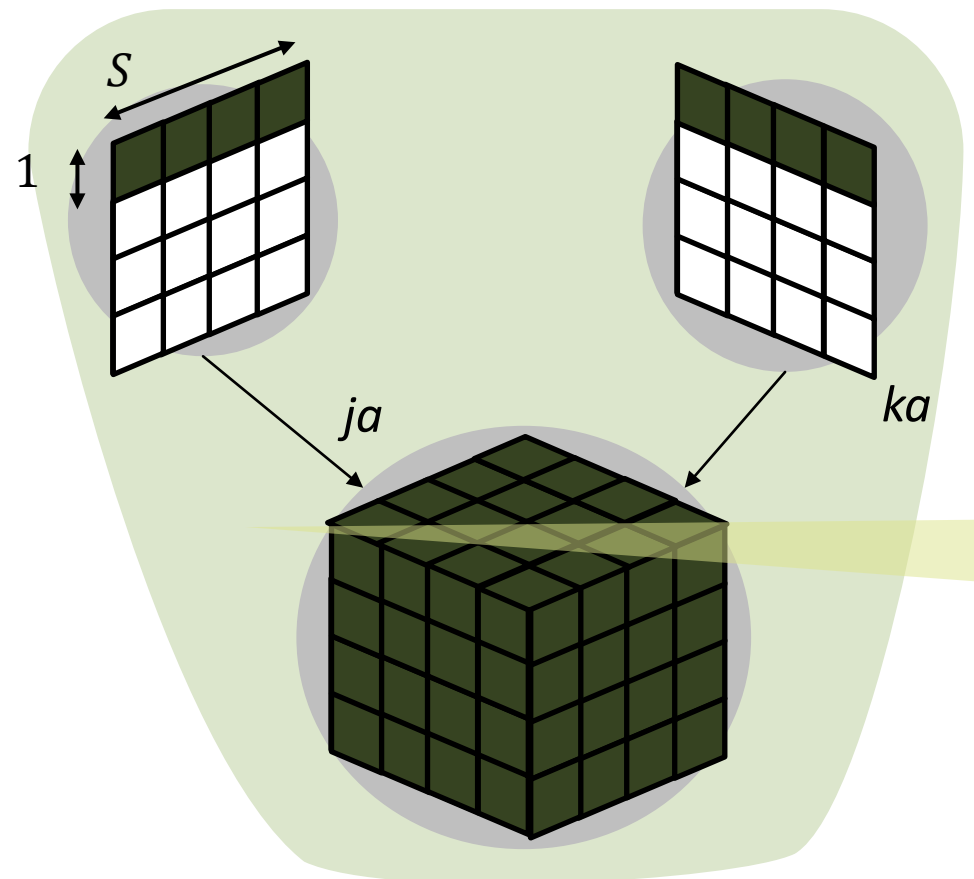$\sqrt{S/3}$

$\sqrt{S/3}$

Inputs are too big to fit in local memory! Can store up to $S$ elements at once!

*Really?* We can do better! We don't reduce **(contract)** mode $a$. No need to keep intermediate results!
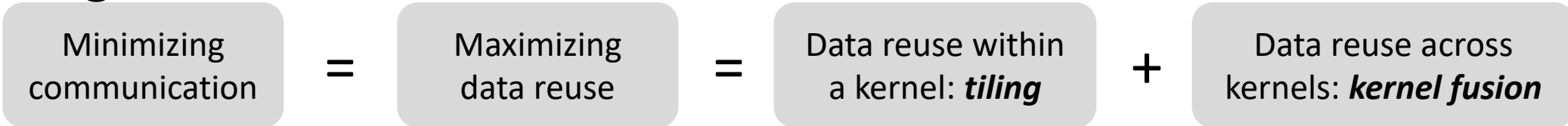
# Minimizing communication ③

| Minimizing communication | = | Maximizing data reuse | = | Data reuse within a kernel: *tiling* | + | Data reuse across kernels: *kernel fusion* |

$$ijk, ja, ka, al \rightarrow il \qquad \boxed{ja, ka \rightarrow jka} \quad ijk, jka \rightarrow ia \qquad ia, al \rightarrow il$$
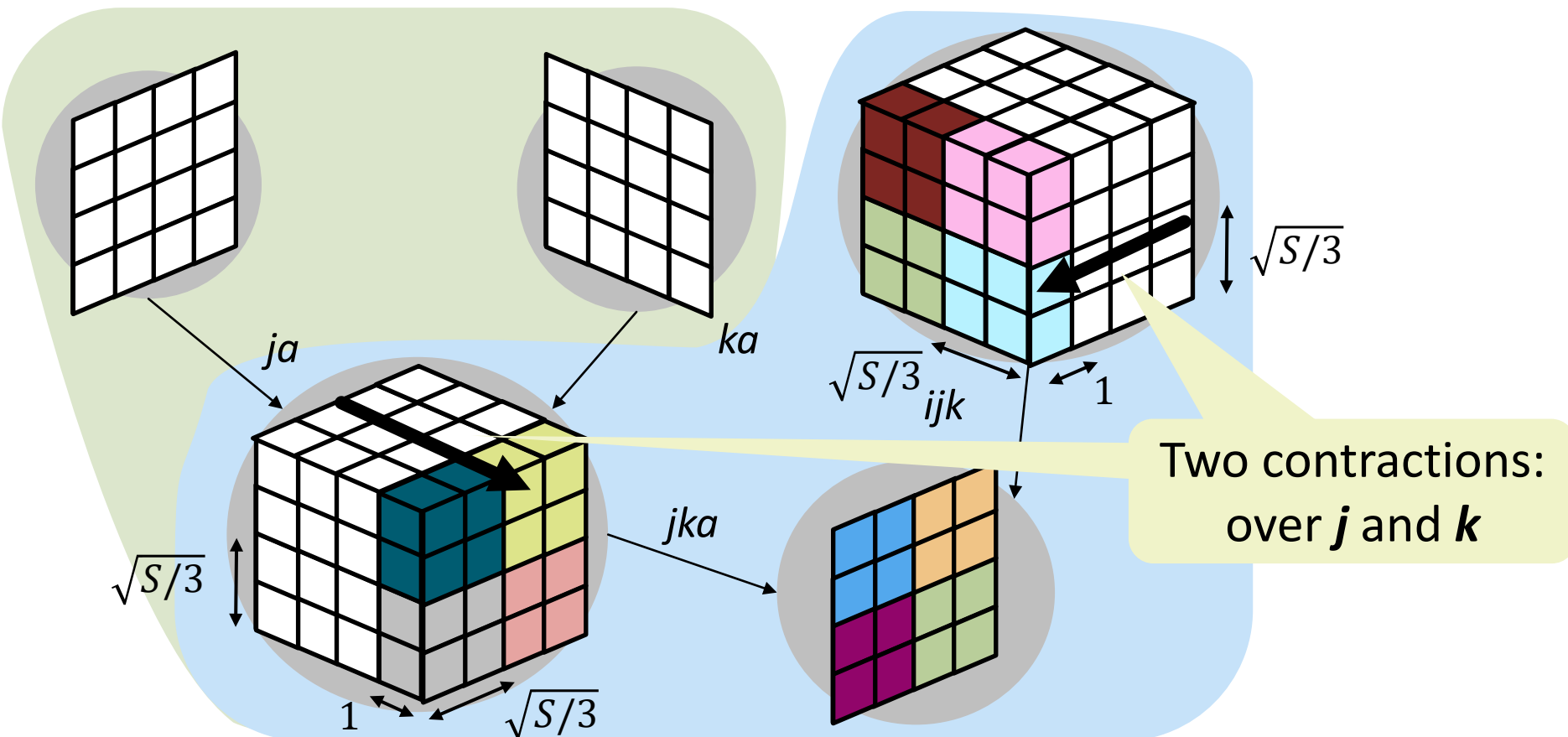


One load operation: $S$ arithmetic operations
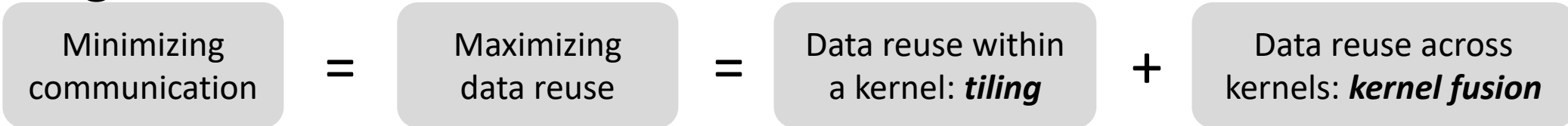
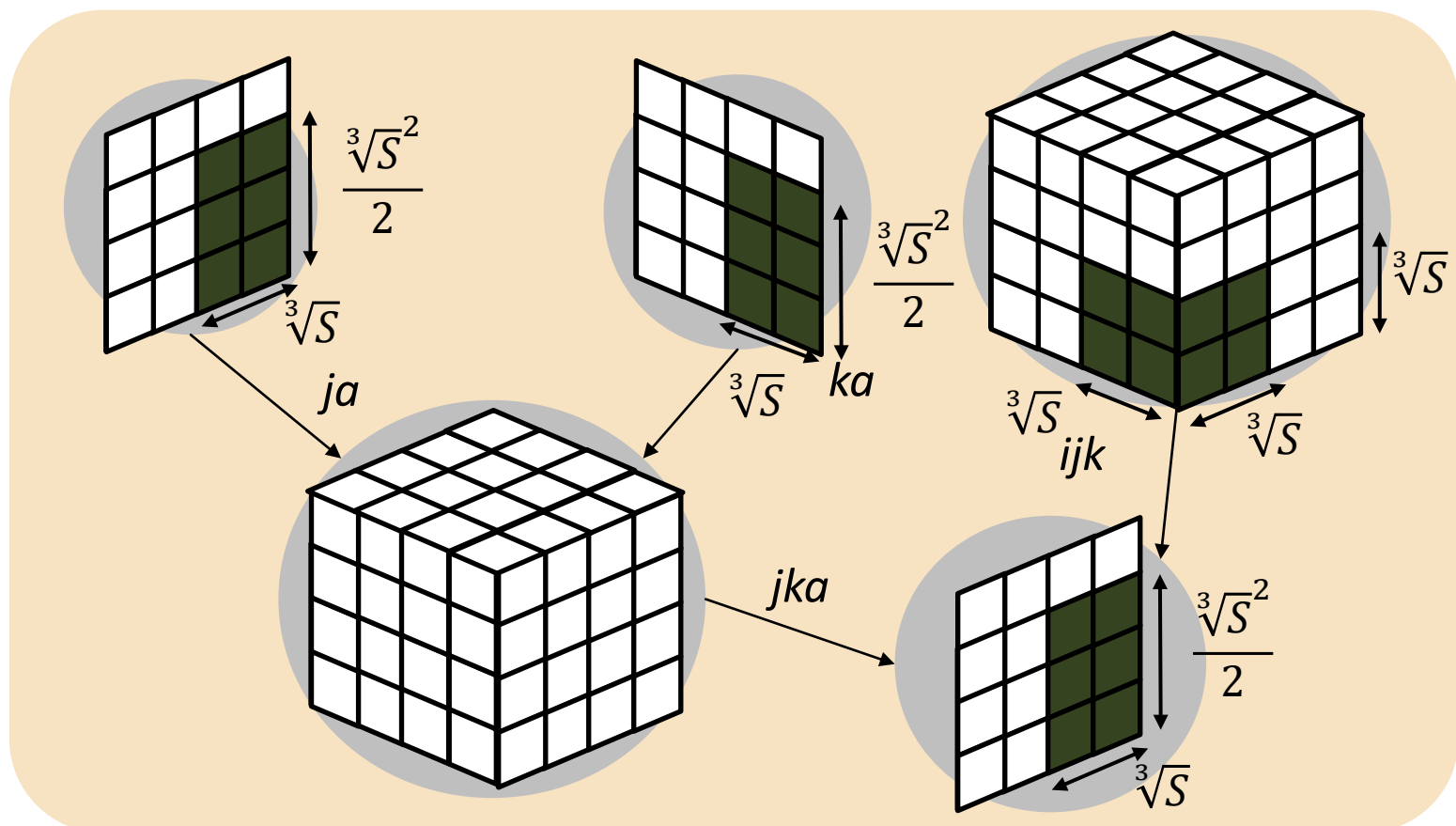(*"cubic" partitioning:* $\sqrt{S}/2$ *arithm. ops per load*)

# Minimizing communication

③

| Minimizing communication | = | Maximizing data reuse | = | Data reuse within a kernel: **tiling** | + | Data reuse across kernels: **kernel fusion** |

$$ijk, ja, ka, al \rightarrow il \qquad ja, ka \rightarrow jka \qquad ijk, jka \rightarrow ia \qquad ia, al \rightarrow il$$



*ja*

*ka*

$\sqrt{S/3}$

*ijk*

$\sqrt{S/3}$

1

$\sqrt{S/3}$

*jka*

$\sqrt{S/3}$

1

$\sqrt{S/3}$

Two contractions: over **j** and **k**

# Minimizing communication

③

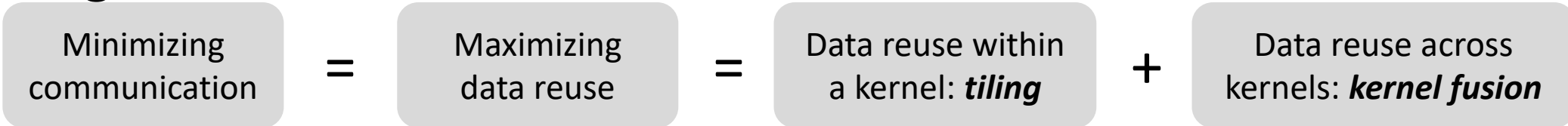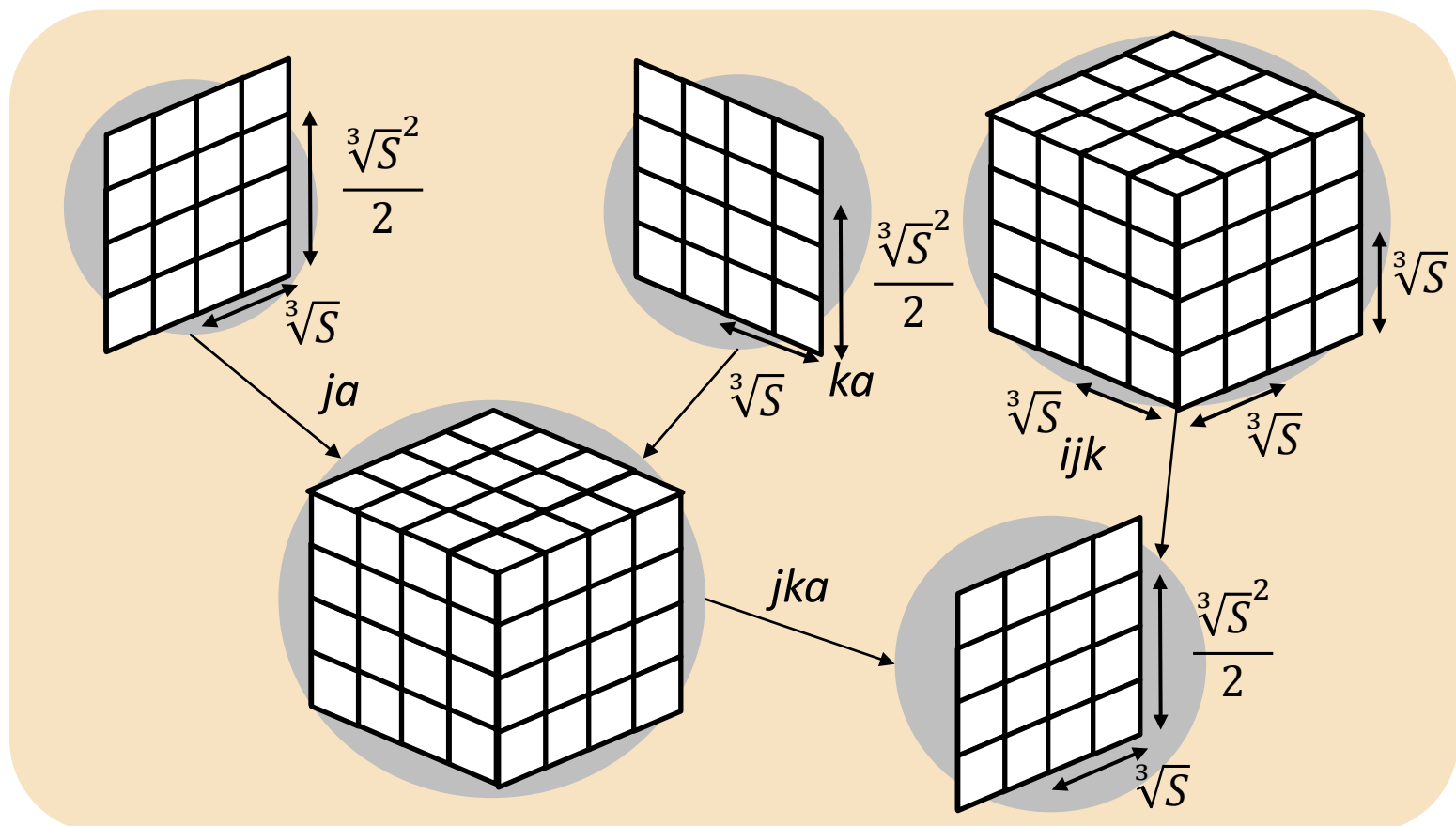| Minimizing communication | = | Maximizing data reuse | = | Data reuse within a kernel: **tiling** | + | Data reuse across kernels: **kernel fusion** |

$$ijk, ja, ka, al \rightarrow il$$   $$ja, ka \rightarrow jka \quad ijk, jka \rightarrow ia$$   $$ia, al \rightarrow il$$
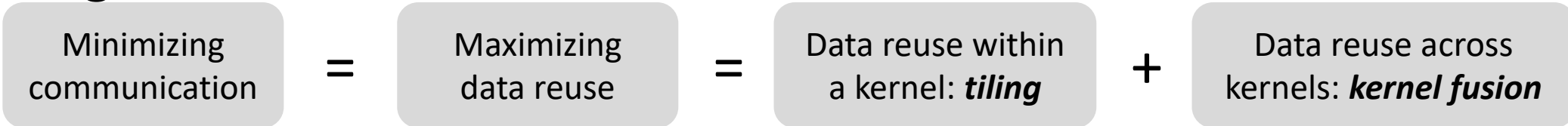
# Minimizing communication ③

| Minimizing communication | = | Maximizing data reuse | = | Data reuse within a kernel: *tiling* | + | Data reuse across kernels: *kernel fusion* |

$$ijk, ja, ka, al \rightarrow il \qquad ja, ka \rightarrow jka \qquad ijk, jka \rightarrow ia \qquad ia, al \rightarrow il$$



**MTTKRP**
*New I/O lower bound:*

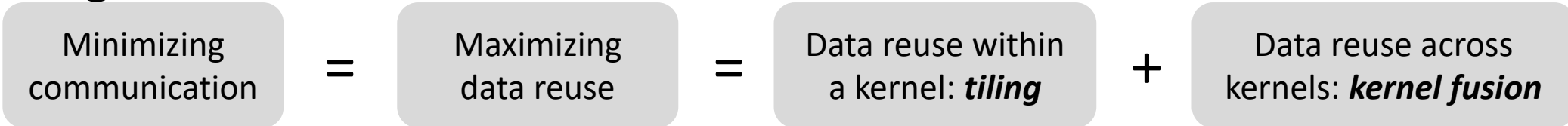$$Q \geq \frac{3 N_1 N_2 N_3 N_4}{S^{2/3}}$$

# Minimizing communication

③

| Minimizing communication | = | Maximizing data reuse | = | Data reuse within a kernel: *tiling* | + | Data reuse across kernels: *kernel fusion* |

$$ijk, ja, ka, al \rightarrow il \qquad ja, ka \rightarrow jka \quad ijk, jka \rightarrow ia \qquad ia, al \rightarrow il$$



$\frac{\sqrt[3]{S}^2}{2}$

$\sqrt[3]{S}$

$ja$

*Asymptotic improvement over "GEMM-ing" tensor contractions.*

$\sqrt[3]{S}$

$jka$

$\frac{\sqrt[3]{S}^2}{2}$

$\sqrt[3]{S}$

**MTTKRP**
*New I/O lower bound:*

$$Q \leq \frac{3N_1 N_2 N_3 N_4}{S^{2/3}}$$

# Minimizing communication

③

| Minimizing communication | = | Maximizing data reuse | = | Data reuse within a kernel: *tiling* | + | Data reuse across kernels: *kernel fusion* |

## Simple Overlap Access Program (SOAP) data movement model

| Iteration vector (3 iteration variables) | $\psi = [\psi^1, \psi^2, \psi^3]$ |
| Access function vector (2 components) | $\phi(\psi) = [\phi_1(\psi), \phi_2(\psi)]$ |
| Iteration variables' ranges | $\psi^1 \in D^1, \psi^2 \in D^2, \psi^3 \in D^3$ |

$$\phi([\mathcal{H}]) = \phi_1([\mathcal{H}]) \cup \phi_2([\mathcal{H}])$$

Iteration domain $\boldsymbol{D} = D^1 \times D^2 \times D^3$

# of points in $\mathcal{H}$ dependent on $\phi_j(\psi_j^*)$
$\theta(\psi_j^*, \mathcal{H})$

"Stretch" in "free" dimensions.
$\max(\theta(\psi_j^*, \mathcal{H})) = \prod_{i \in \Psi_j'} |D^i|$

"Stretch" in remaining dimensions. $\delta_j^{-1}$ is constant.

"Almost rectangular"

Extend in remaining domains

$\phi_j[\mathcal{H}_r] = \phi_j[\mathcal{H}]$
$\delta_j^{-1}(\mathcal{H}_r) \leq \delta_j^{-1}(\mathcal{H})$

$\delta_j^{-1}(\mathcal{H}_{rec}) = \frac{|\phi_j[\mathcal{H}_{rec}]|}{|\mathcal{H}_{rec}|} = \frac{|\phi_j[\mathcal{H}_r]|}{|\mathcal{H}_r|}$
$\delta_j^{-1}(\mathcal{H}_{rec}) = \delta_j^{-1}(\mathcal{H}_r) \leq \delta_j^{-1}(\mathcal{H})$

$H = \{C, E\}$

A   B   D

$\phi_{St1,1}$   $\phi_{St1,2}$

Executing statements **St₁** and **St₂** "together"

C

$\phi_{St2,3}$

$\phi_{St2,2}$

Array **C** is not an input anymore!

E   $\phi_{St2,1}$

| MTTKRP definition | $u_{il} = \sum_{j,k} t_{ijk} v_{jl} w_{kl}$ |
| `opt_einsum` decomposition | 1.  $x_{jkl} = v_{jl} w_{kl}$ <br> 2.  $u_{il} = \sum_{j,k} t_{ijk} x_{jkl}$ |

SDG

v   w
$\mathcal{V} = JL$   $\mathcal{W} = KL$

t   x
$\mathcal{T} = IJK$   $\mathcal{X} = JKL$

u

I/O lower bound

$\max I \cdot J \cdot K \cdot L \quad \text{s.t.}$
$I \cdot J \cdot K + J \cdot L + K \cdot L \leq X$

Solution:
$\rho = \frac{S^{2/3}}{3}$
$Q = \frac{|V|}{\rho} = \frac{3N_1 N_2 N_3 N_4}{S^{2/3}}$

## MTTKRP
*New I/O lower bound:*
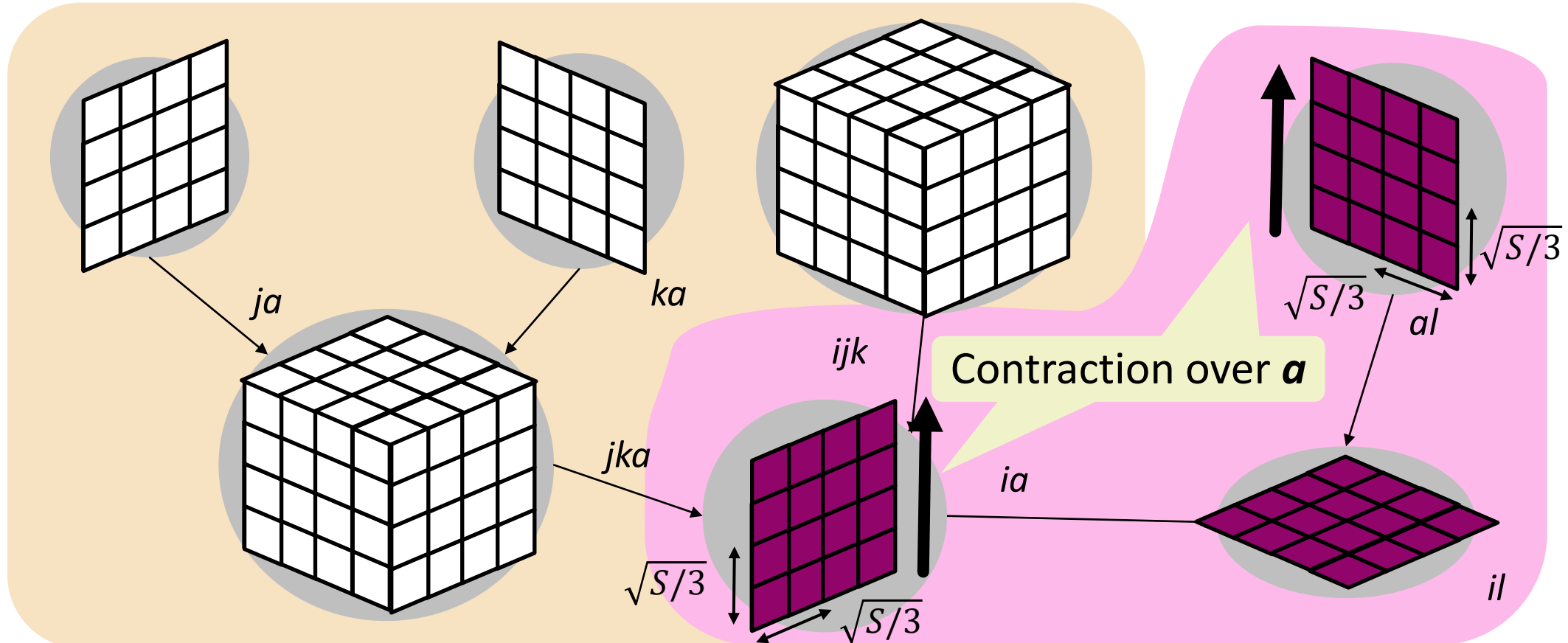
$$Q \geq \frac{3N_1 N_2 N_3 N_4}{S^{2/3}}$$

# Minimizing communication

③

| Minimizing communication | = | Maximizing data reuse | = | Data reuse within a kernel: **tiling** | + | Data reuse across kernels: **kernel fusion** |

$$ijk, ja, ka, al \rightarrow il \qquad ja, ka \rightarrow jka \qquad ijk, jka \rightarrow ia \qquad ia, al \rightarrow il$$



*ja*

*ka*

*ijk*

Contraction over **a**

*jka*

*ia*

*il*

$\sqrt{S/3}$

$\sqrt{S/3}$

$\sqrt{S/3}$

$\sqrt{S/3}$

*al*

20

# Minimizing communication

③

| Minimizing communication | = | Maximizing data reuse | = | Data reuse within a kernel: *tiling* | + | Data reuse across kernels: *kernel fusion* |

## COMMUNICATION OPTIMAL:

Tiling

Kernel fusion

Parallel data decomposition

## For a given input einsum

# Deinsum: Practically I/O Optimal Multi-Linear Algebra

**Simple Overlap Access Pattern (SOAP)**
**data movement model**

① **Input**

$$ijk, ja, ka, al \rightarrow il$$

einsum, string

② **Split to binary operations**

$$ja, ka \rightarrow jka$$
$$ijk, jka \rightarrow ia$$
$$ia, al \rightarrow il$$

uses operation associativity to minimize #ops

③ **Communication-optimal schedules**



*ja*   *ka*   *ijk*

*jka*   *ia*   *al*

*il*

**Optimal tiling, parallel distribution, and kernel fusion**

# Deinsum: Practically I/O Optimal Multi-Linear Algebra

**block-distributions**

④ **Iteration spaces and distribution**

Iteration space partition:
`MPI_Cart_sub`

Distribute initial data:
`MPI_Broadcast`

⑤ **Automated code generation**

```
for j in range(NJ//P0J):
  for k in range(NK//P0K):
    for a in range(NA//P0A):
      t0[j, k, a] += A[j, a] * B[k, a]

t1 = np.tensordot(X, t0,
        axes=([1, 2], [0, 1])
mpi.Allreduce(t1, comm=grid0_t1)

t2 = deinsum.Redistribute(
        comm1=grid0, comm2=grid1)
```

**Auto generating DaCe-Python code compiled to shared library**

⑥ **Results**

- **up to 19x speedup over SotA**
- **CPU and GPU support**

# Grouping Operations

③➡④

```python
for j in range(NJ):
    for k in range(NK):
        for a in range(NA):
            t0[j,k,a]+=A[j,a]*B[k,a]
```

```python
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

```python
out = t1 @ C
```

# Iteration Spaces: Global View

④

Cartesian process grid

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

```
for j in range(NJ):
    for k in range(NK):
        for a in range(NA):
            t0[j,k,a]+=A[j,a]*B[k,a]
```

```
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

# Iteration Spaces: Global View

④

**Cartesian process grid**

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

```
for j in range(p0j*NJ//P0J,(p0j+1)*NJ//P0J):
  for k in range(p0k*NK//P0K,(p0k+1)*NJ//P0J):
    for a in range(p0a*NA//P0A,(p0a+1)*NA//P0A):
      t0[j,k,a]+=A[j,a]*B[k,a]
```

```
t1[p0i*NI//P0I,(p0i+1)*NI//P0I,
   p0a*NA//P0A,(p0a+1)*NA//P0A] = np.tensordot(
   X[p0i*NI//P0I,(p0i+1)*NI//P0I,
     p0j*NJ//P0J,(p0j+1)*NJ//P0J,
     p0k*NK//P0K,(p0k+1)*NJ//P0J],
   t0[p0j*NJ//P0J,(p0j+1)*NJ//P0J,
      p0k*NK//P0K,(p0k+1)*NJ//P0J,
      p0a*NA//P0A,(p0a+1)*NA//P0A],
   axes=([1,2], [0,1]))
```

$P_a^{(0)}$

$P_i^{(0)}$

$p_k^{(0)}$

$P_j^{(0)}$

$p_j^{(0)}$

$P_k^{(0)}$

$p_i^{(0)}$

$p_a^{(0)}$

**process-local slices – global coordinates**

$N_i/P_i^{(0)}$

$N_j/P_j^{(0)}$

$N_k/P_k^{(0)}$

# Iteration Spaces: Local View

④

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$
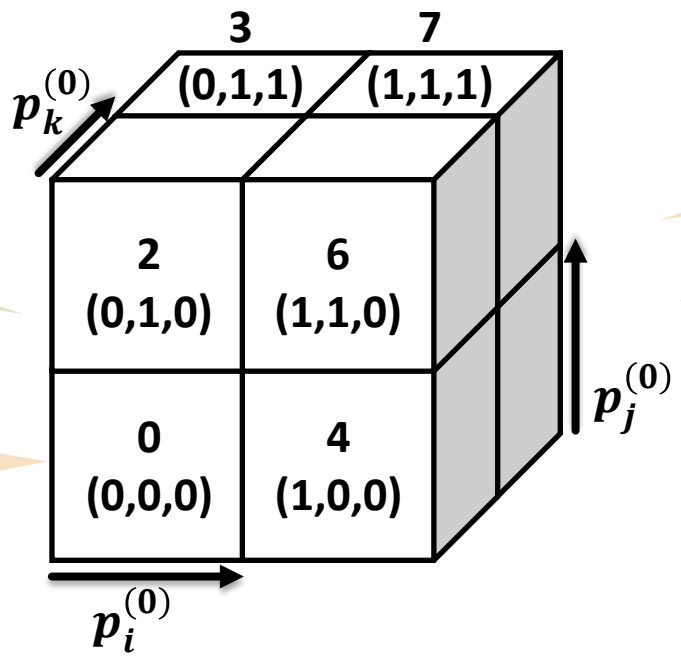
```python
for j in range(NJ//P0J):
    for k in range(NK//P0K):
        for a in range(NA//P0A):
            t0[j,k,a]+=A[j,a]*B[k,a]
```

```python
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

# Iteration Spaces: Local View

process-local slices – local coordinates

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$
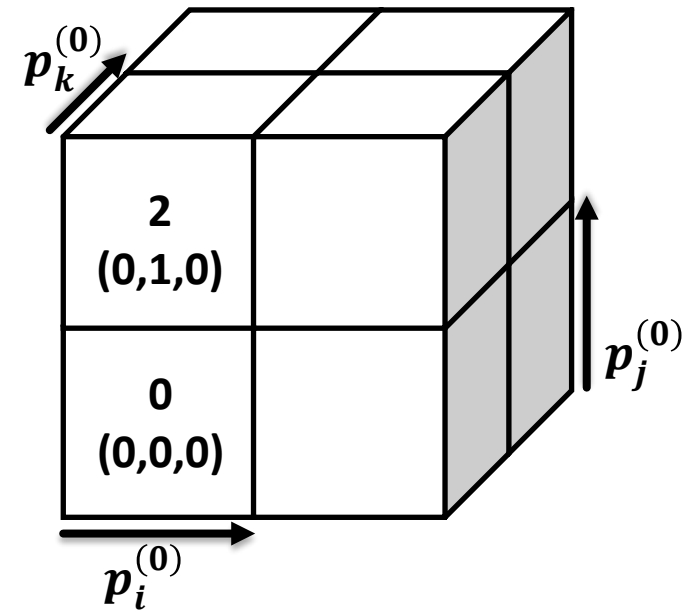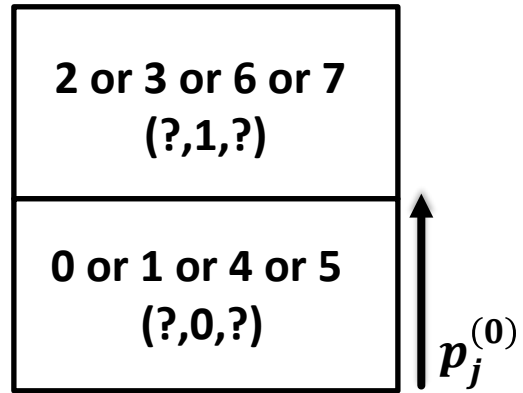
```python
for j in range(NJ//P0J):
    for k in range(NK//P0K):
        for a in range(NA//P0A):
            t0[j,k,a]+=A[j,a]*B[k,a]
```

```python
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

```c
//                  {  i,   j,   k,   a}
int dims[4] =      {P0I, P0J, P0K, P0A};
int periods[4] = {  0,   0,   0,   0};
MPI_Comm grid0;
MPI_Cart_create(MPI_COMM_WORLD, 4, dims,
                 periods, false, &grid0);
```

# Iteration Spaces: Local View

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

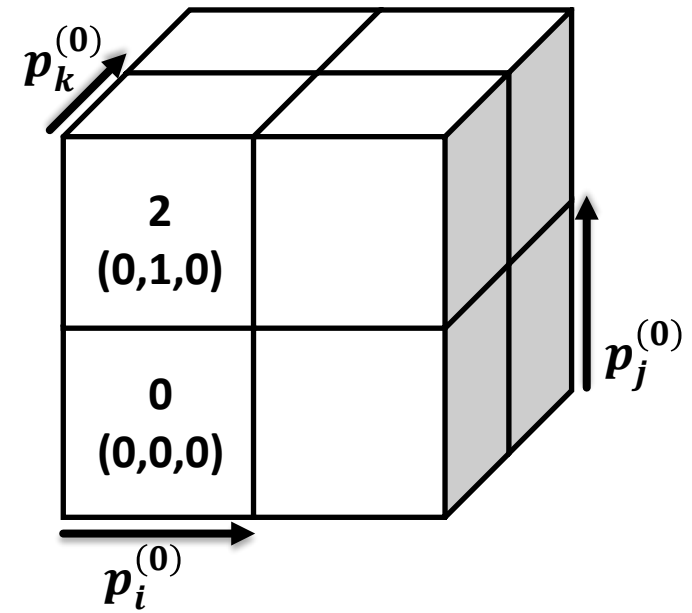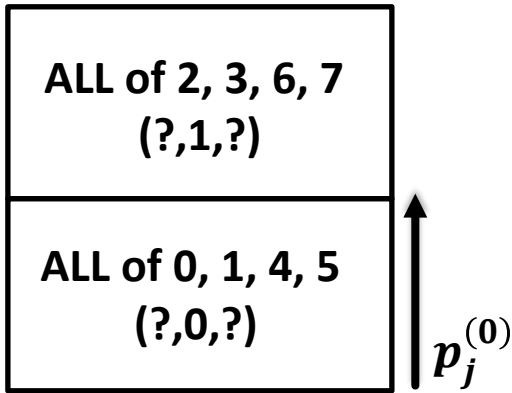process-local slices – local coordinates

```python
for j in range(NJ//P0J):
    for k in range(NK//P0K):
        for a in range(NA//P0A):
            t0[j,k,a]+=A[j,a]*B[k,a]

t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

```c
//                   {  i,   j,   k,   a}
int dims[4] =    {P0I, P0J, P0K, P0A};
int periods[4] = {  0,   0,   0,   0};
MPI_Comm grid0;
MPI_Cart_create(MPI_COMM_WORLD, 4, dims,
                  periods, false, &grid0);
```

$$P = P_i^{(1)} P_l^{(1)} P_a^{(1)}$$

```python
out = t1 @ C
```

```c
//                   {  i,   l,   a}
int dims[3] =    {P1I, P1L, P1A};
int periods[3] = {  0,   0,   0};
MPI_Comm grid1;
MPI_Cart_create(MPI_COMM_WORLD, 3, dims,
                  periods, false, &grid1);
```

# Iteration Spaces: Practically I/O Optimal Distribution

④

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

```python
for j in range(NJ//P0J):
    for k in range(NK//P0K):
        for a in range(NA//P0A):
            t0[j,k,a]+=A[j,a]*B[k,a]
```
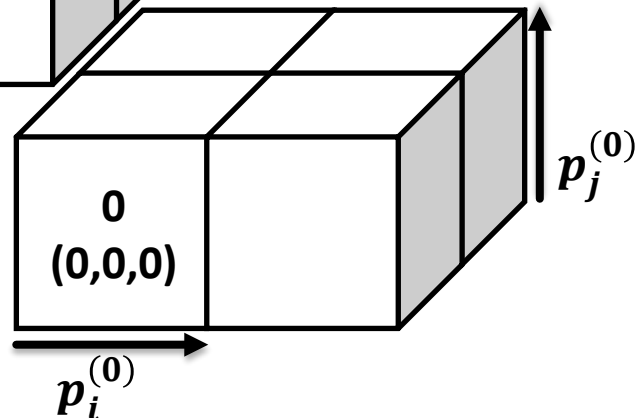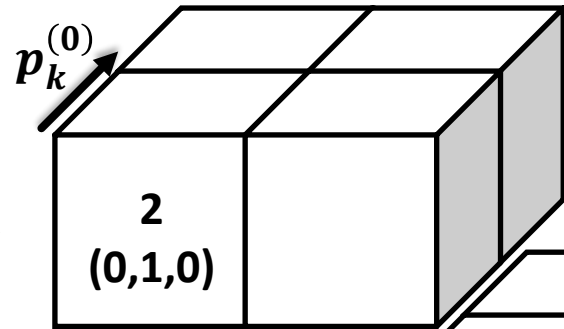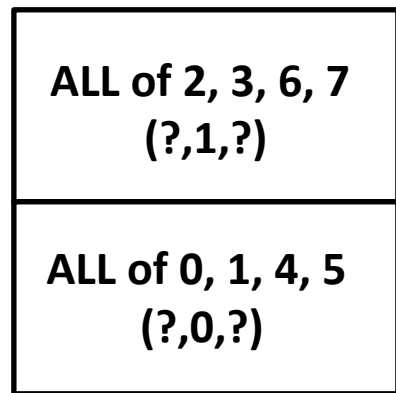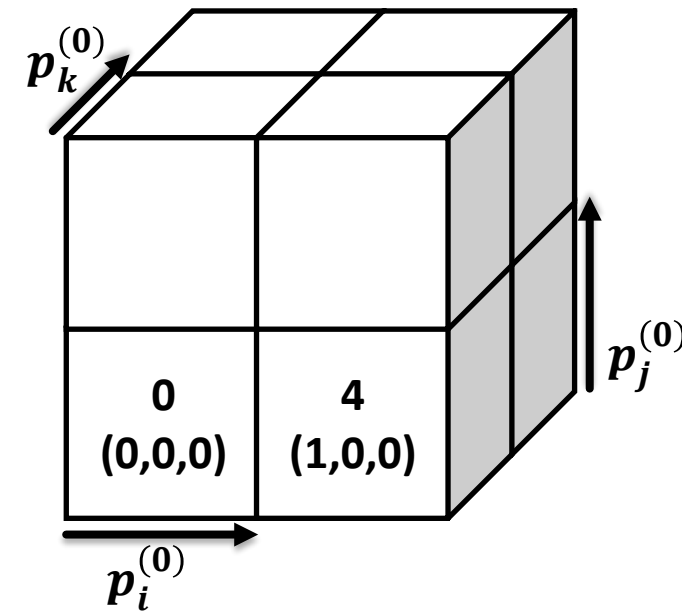
```python
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

**Optimal Tile Sizes**

$$T_i = T_j = T_k = S^{1/3}, T_a = \frac{S^{2/3}}{2}$$

```python
grid_ijka = {
    #       [i, j, k, a]
    1:      [1, 1, 1, 1],
    2:      [1, 1, 2, 1],
    4:      [1, 2, 2, 1],
    8:      [2, 2, 2, 1],
    12:     [2, 2, 3, 1],
    16:     [2, 2, 4, 1],
    27:     [3, 3, 3, 1],
    32:     [2, 4, 4, 1],
    64:     [4, 4, 4, 1],
    125:    [5, 5, 5, 1],
    128:    [4, 4, 8, 1],
    252:    [6, 6, 7, 1],
    256:    [4, 8, 8, 1],
    512:    [8, 8, 8, 1],
}
```
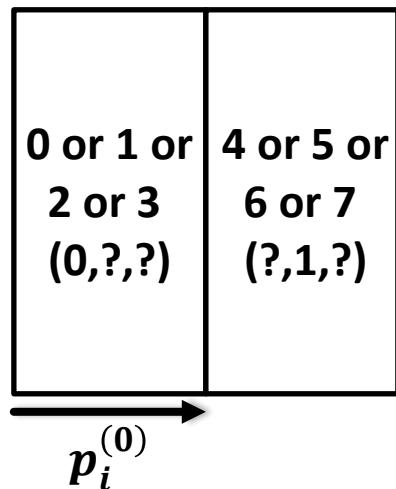
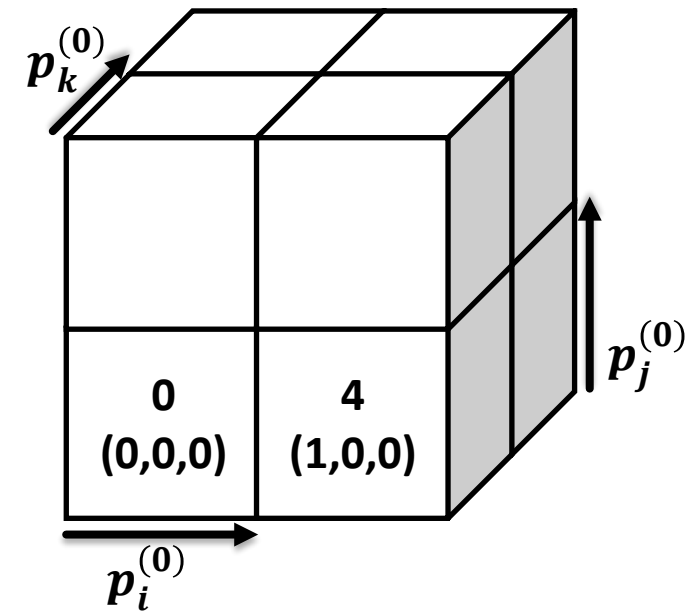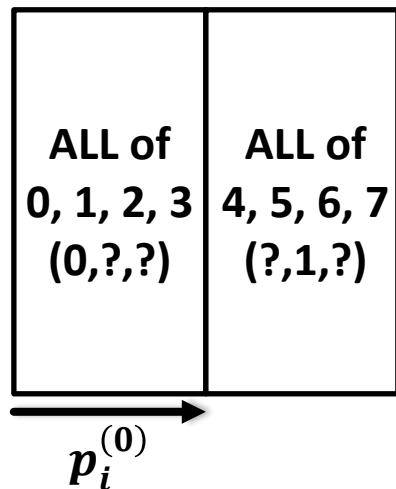# Data Distribution: Replication

④

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

```
for j in range(NJ//P0J):
    for k in range(NK//P0K):
        for a in range(NA//P0A):
            t0[j,k,a]+=A[j,a]*B[k,a]
```

```
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

A[j,a]

| |
|---|
| 2 or 3 or 6 or 7 (?,1,?) |
| 0 or 1 or 4 or 5 (?,0,?) |

$p_j^{(0)}$

$p_k^{(0)}$

2 (0,1,0)

0 (0,0,0)

$p_j^{(0)}$

$p_i^{(0)}$

# Data Distribution: Replication

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$
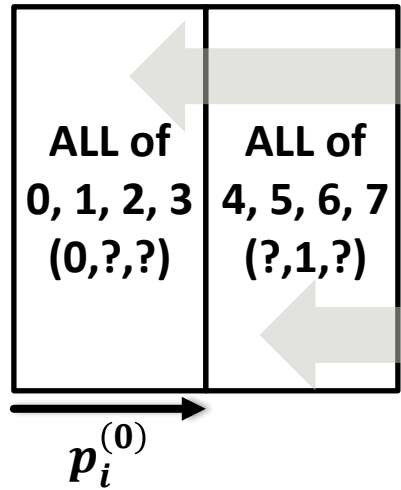
```python
for j in range(NJ//P0J):
    for k in range(NK//P0K):
        for a in range(NA//P0A):
            t0[j,k,a]+=A[j,a]*B[k,a]
```

```python
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

A[j,a]

| |
|---|
| **ALL of 2, 3, 6, 7** <br> **(?,1,?)** |
| **ALL of 0, 1, 4, 5** <br> **(?,0,?)** |

$p_j^{(0)}$

$p_k^{(0)}$

**2** <br> **(0,1,0)**

**0** <br> **(0,0,0)**

$p_j^{(0)}$

$p_i^{(0)}$

# Data Distribution: Replication

④

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

```python
for j in range(NJ//P0J):
    for k in range(NK//P0K):
        for a in range(NA//P0A):
            t0[j,k,a]+=A[j,a]*B[k,a]
```

```python
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

```c
//                  {i, j, k, a}
int remain_A[4] = {1, 0, 1, 0};
MPI_Comm grid0_A;
MPI_Cart_sub(grid0, remain_A, &grid0_A);
MPI_Bcast(A, count, datatype, leader,
          grid0_A);
```

A[j,a]



$p_k^{(0)}$

ALL of 2, 3, 6, 7
(?,1,?)

replication

2
(0,1,0)

ALL of 0, 1, 4, 5
(?,0,?)

replication

0
(0,0,0)

$p_j^{(0)}$

$p_j^{(0)}$

$p_i^{(0)}$

# Data Distribution: Partial Sums

④

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

```
for j in range(NJ//P0J):
    for k in range(NK//P0K):
        for a in range(NA//P0A):
            t0[j,k,a]+=A[j,a]*B[k,a]
```

```
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

t1[i,a]

| 0 or 1 or 2 or 3 (0,?,?) | 4 or 5 or 6 or 7 (?,1,?) |
|---|---|

$p_i^{(0)}$



$p_k^{(0)}$

$p_j^{(0)}$

| 0 (0,0,0) | 4 (1,0,0) |
|---|---|

$p_i^{(0)}$

# Data Distribution: Partial Sums

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

```
for j in range(NJ//P0J):
    for k in range(NK//P0K):
        for a in range(NA//P0A):
            t0[j,k,a]+=A[j,a]*B[k,a]
```

```
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

t1[i,a]

output has partial sums

| ALL of 0, 1, 2, 3 (0,?,?) | ALL of 4, 5, 6, 7 (?,1,?) |
|---|---|

$p_i^{(0)}$

$p_k^{(0)}$

$p_j^{(0)}$

| 0 (0,0,0) | 4 (1,0,0) |
|---|---|

$p_i^{(0)}$

# Data Distribution: Partial Sums

④

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

```python
for j in range(NJ//P0J):
    for k in range(NK//P0K):
        for a in range(NA//P0A):
            t0[j,k,a]+=A[j,a]*B[k,a]
```

```python
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

```c
//                    {i, j, k, a}
int remain_t1[4] = {0, 1, 1, 0};
MPI_Comm grid0_t1;
MPI_Cart_sub(grid0, remain_t1, &grid0_t1);
MPI_Allreduce(MPI_IN_PLACE, t1, count,
              datatype, leader, MPI_SUM,
              grid0_t1);
```

`t1[i,a]`

**output has partial sums**

**reduction**

| ALL of 0, 1, 2, 3 (0,?,?) | ALL of 4, 5, 6, 7 (?,1,?) |

**reduction**

$p_i^{(0)}$

$p_k^{(0)}$

0
(0,0,0)

$p_i^{(0)}$

4
(1,0,0)

$p_j^{(0)}$

# Data Redistribution

④

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

```
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

# Data Redistribution

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

```
t1 = np.tensordot(X, t0, axes=([1,2], [0,1]))
```

t1[i,a]

$p_k^{(0)}$

reduction

| ALL of 0, 1, 2, 3 (0,?,?) | ALL of 4, 5, 6, 7 (?,1,?) |

2 (0,1,

reduction

0 (0,0,

$p_i^{(0)}$

$p_i^{(0)}$

$$P = P_i^{(1)} P_l^{(1)} P_a^{(1)}$$

```
out = t1 @ C
```

```
//
int dims[3
int period
MPI_Comm g
MPI_Cart_cre
                                    }
                        periods, false, &grid1);
```

```
grid_ila = {
    #       [i, l, a]
    1:      [1, 1, 1],
    2:      [1, 1, 2],
    4:      [1, 2, 2],
    8:      [2, 2, 2],
    12:     [2, 2, 3],
    16:     [2, 2, 4],
    27:     [3, 3, 3],
    32:     [2, 4, 4],
    64:     [4, 4, 4],
    125:    [5, 5, 5],
    128:    [4, 4, 8],
    252:    [6, 6, 7],
    256:    [4, 8, 8],
    512:    [8, 8, 8],
```
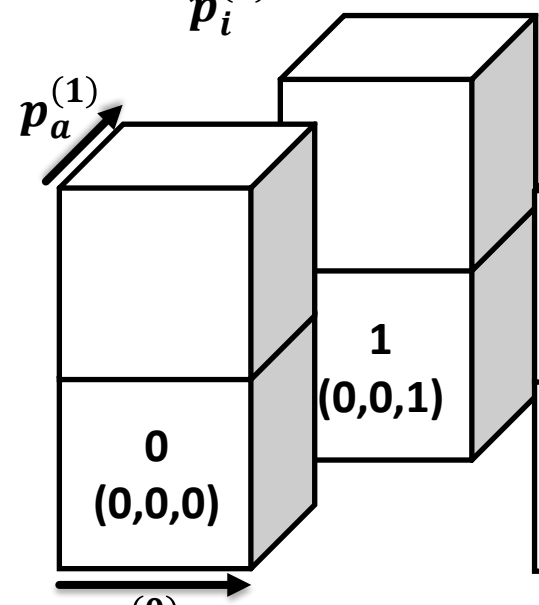
# Data Redistribution

```
grid_ijka = {
    #      [i, j, k, a]
    1:     [1, 1, 1, 1],
    2:     [1, 1, 2, 1],
    4:     [1, 2, 2, 1],
    8:     [2, 2, 2, 1],
    12:    [2, 2, 3, 1],
    16:    [2, 2, 4, 1],
    27:    [3, 3, 3, 1],
    32:    [2, 4, 4, 1],
    64:    [4, 4, 4, 1],
    125:   [5, 5, 5, 1],
    128:   [4, 4, 8, 1],
    252:   [6, 6, 7, 1],
    256:   [4, 8, 8, 1],
    512:   [8, 8, 8, 1],
}
```

```
grid_ila = {
    #      [i, l, a]
    1:     [1, 1, 1],
    2:     [1, 1, 2],
    4:     [1, 2, 2],
    8:     [2, 2, 2],
    12:    [2, 2, 3],
    16:    [2, 2, 4],
    27:    [3, 3, 3],
    32:    [2, 4, 4],
    64:    [4, 4, 4],
    125:   [5, 5, 5],
    128:   [4, 4, 8],
    252:   [6, 6, 7],
    256:   [4, 8, 8],
    512:   [8, 8, 8],
}
```

$t1 =$     ,1]))

$p_k^{(0)}$

2
(0,1,

0
(0,0,

tion

$p_i^{(0)}$

```
//
int dims[3
int period
MPI_Comm g
MPI_Cart_cr                      ,
                 periods, false, &grid1);
```

# Data Redistribution

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

| | |
|---|---|
| 0, 1, 2, 3 (0,?,?) | 4, 5, 6, 7 (?,1,?) |

$p_i^{(0)}$

t1[i,a]

$$P = P_i^{(1)} P_l^{(1)} P_a^{(1)}$$

| | |
|---|---|
| 1, 3 (0,?,0) | 5, 7 (0,?,0) |
| 0, 2 (0,?,0) | 4, 6 (0,?,0) |

$p_a^{(1)}$ $p_i^{(1)}$



$p_k^{(0)}$

0 (0,0,0)

4 (1,0,0)

$p_j^{(0)}$

$p_i^{(0)}$

$p_a^{(1)}$

1 (0,0,1)

5 (1,0,1)

0 (0,0,0)

4 (1,0,0)

$p_l^{(1)}$

$p_i^{(0)}$

# Data Redistribution

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

| | |
|---|---|
| 0, 1, 2, 3 (0,?,?) | 4, 5, 6, 7 (?,1,?) |

$p_i^{(0)}$

sub-communicator root processes send/recv blocks

$p_k^{(0)}$

0 (0,0,0)

$p_i^{(0)}$

4 (1,0,0)

$p_j^{(0)}$

t1[i,a]

$$P = P_i^{(1)} P_l^{(1)} P_a^{(1)}$$

| | |
|---|---|
| 1, 3 (0,?,0) | 5, 7 (0,?,0) |
| 0, 2 (0,?,0) | 4, 6 (0,?,0) |

$p_a^{(1)}$

$p_i^{(1)}$

$p_a^{(1)}$

2 (0,1,0)

1 (0,0,1)

0 (0,0,0)

4 (1,0,0)

5 (1,0,1)

$p_l^{(1)}$

$p_i^{(0)}$

# Data Redistribution

$$P = P_i^{(0)} P_j^{(0)} P_k^{(0)} P_a^{(0)}$$

| 0, 1, 2, 3 (0,?,?) | 4, 5, 6, 7 (?,1,?) |
|---|---|

$p_i^{(0)}$

t1[i,a]

$$P = P_i^{(1)} P_l^{(1)} P_a^{(1)}$$

| 1, 3 (0,?,0) | 5, 7 (0,?,0) |
|---|---|
| 0, 2 (0,?,0) | 4, 6 (0,?,0) |

$p_a^{(1)}$   $p_i^{(1)}$

sub-communicator root processes send/recv blocks

$p_k^{(0)}$

0 (0,0,0)

$p_i^{(0)}$

4 (1,0,0)

$p_j^{(0)}$

and broadcast to the rest

$p_a^{(1)}$

0 (0,0,0)

1 (0,0,1)

4 (1,0,0)

5 (1,0,1)

$p_l^{(1)}$

$p_i^{(0)}$

# Automated code generation: from DaCe-Python to C++

⑤

```python
grid0 = mpi.Cart_create(dims=[P0I,P0J,P0K,P0A])
grid0_t1 = mpi.Cart_sub(comm=grid0, remain=[False,True,True,False])
grid1 = mpi.Cart_create(dims=[P1I, P1L, P1A])
grid1_out = mpi.Cart_sub(comm=grid1, remain=[False,False,True])
```

**MPI comm interface**

```python
# ja,ka->jka
t0 = np.zeros((NJ//P0J, NK//P0K, NA//P0A), dtype=X.dtype)
for j in range(NJ//P0J):
  for k in range(NK//P0K):
    for a in range(NA//P0A):
      t0[j, k, a] += A[j, a] * B[k, a]
```

**products without contraction: for-loops**

```python
# ijk,jka->ia
t1 = np.tensordot(X, t0, axes=([1, 2], [0, 1]))
mpi.Allreduce(t1, comm=grid0_t1, op=mpi.SUM)
```

**dot products: tensordot**

**MPI collectives**

```python
t2 = deinsum.Redistribute(t1, comm1=grid0, comm2=grid1)
```

**data redistribution**

```python
# ia,al->il
out = t2 @ C
mpi.Allreduce(out, comm=grid1_out)
```

49

⑤

```cpp
int grid1_remain[4] = {0, 1, 1, 0};
MPI_Cart_sub(grid0_comm, pgrid1_remain, &grid1_comm);
MPI_Comm_rank(grid1_comm, &grid1_rank);
MPI_Cart_coords(grid1_comm, grid_1_rank, 2, grid_1_coords);
```

**(CUDA-aware) MPI**

```cpp
#pragma omp parallel for
for (auto j = 0; j < S1; j += 1)
  for (auto k = 0; k < S2; k += 1)
    for (auto a = 0; a < S3; a += 1)
      t0[S3*S2*j + S3*k + a] = A[S3*j + a]* B[S3*k + a];
```

**CPU – OpenMP**
**GPU – CUDA kernels**

```cpp
cblas_dgemm(CblasColMajor, CblasNoTrans, CblasNoTrans,
            S3, S0, S1*S2,
            1.0, t0, S3, X, S1*S2,
            0.0, t1, S3);
```

**CPU – TTGT**
**GPU – cuTENSOR**

```cpp
MPI_Allreduce(MPI_IN_PLACE, out, S0*S3, MPI_DOUBLE, MPI_SUM,
              grid1_comm);
```

# Results: CPU

weak scaling
Piz Daint 1 – 512 nodes
Intel E5-2690 v3 (12 cores)
Deinsum
CTF

total runtime

compute runtime

**3MM**

2.4x

**MTTKRP-O5-M0**

18.1x

**TTMc-O5-M0**

16.0x

Runtime [s]

Number of Nodes

# Results: CPU

weak scaling
Piz Daint 1 – 512 nodes
Intel E5-2690 v3 (12 cores)
Deinsum
CTF

total runtime

compute runtime

bottleneck: collectives on sub-communicators

2.4x

3MM

18.1x

MTTKRP-O5-M0

16.0x

TTMc-O5-M0

Runtime [s]

64 tiles

"tiling" dimensions

1 tile

512:  [8, 8, 8]

512:  [1, 16, 32, 1, 1, 1]

# Results: GPU

**runtime with data already in GPU global memory**

**total runtime**

weak scaling
Piz Daint 1 – 512 nodes
Nvidia P100
Deinsum
CTF

**3MM**

**2.5x**

**MTTKRP-O5-M0**

**9.0x**

**TTMc-O5-M0**

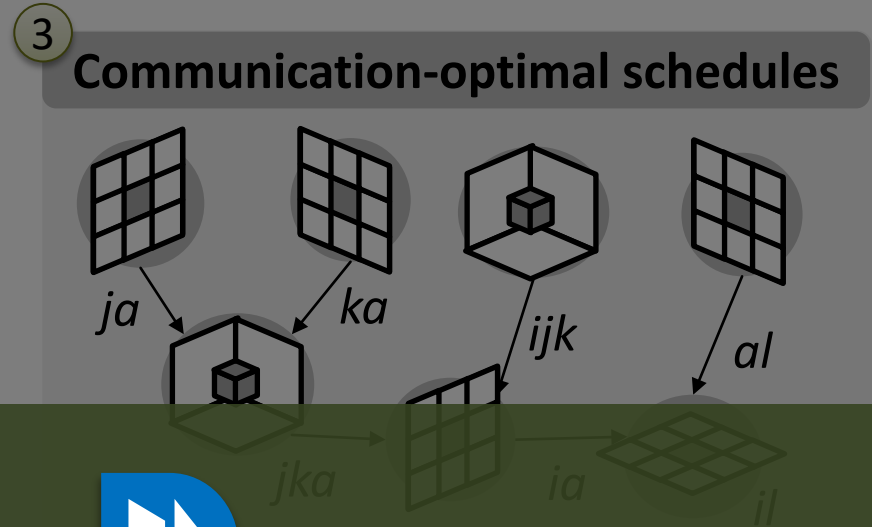**3.9x**

Runtime [s]

Number of Nodes

# Conclusions

① **Input**

$$ijk, ja, ka, al \rightarrow il$$

② **Split to binary operations**

$$ja, ka \rightarrow jka$$
$$ijk, jka \rightarrow ia$$
$$ia, al \rightarrow il$$

③ **Communication-optimal schedules**



**github.com/spcl/dace**

④ Iteration spaces and distribution

Iteration space partition:
`MPI_Cart_sub`

Distribute initial data:
`MPI_Broadcast`

⑤ Automated code generation

```
for j in range(NJ//P0J):
  for k in range(NK//P0K):
    for a in range(NA//P0A):
      t0[j, k, a] += A[j, a] * B[k, a]

t1 = np.tensordot(X, t0,
          axes=([1, 2], [0, 1])
mpi.Allreduce(t1, comm=grid0_t1)

t2 = deinsum.Redistribute(t1,
          comm1=grid0, comm2=grid1)
```

⑥ Results

- **up to 19x speedup over CTF**
- **CPU and GPU support**