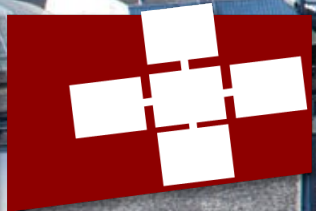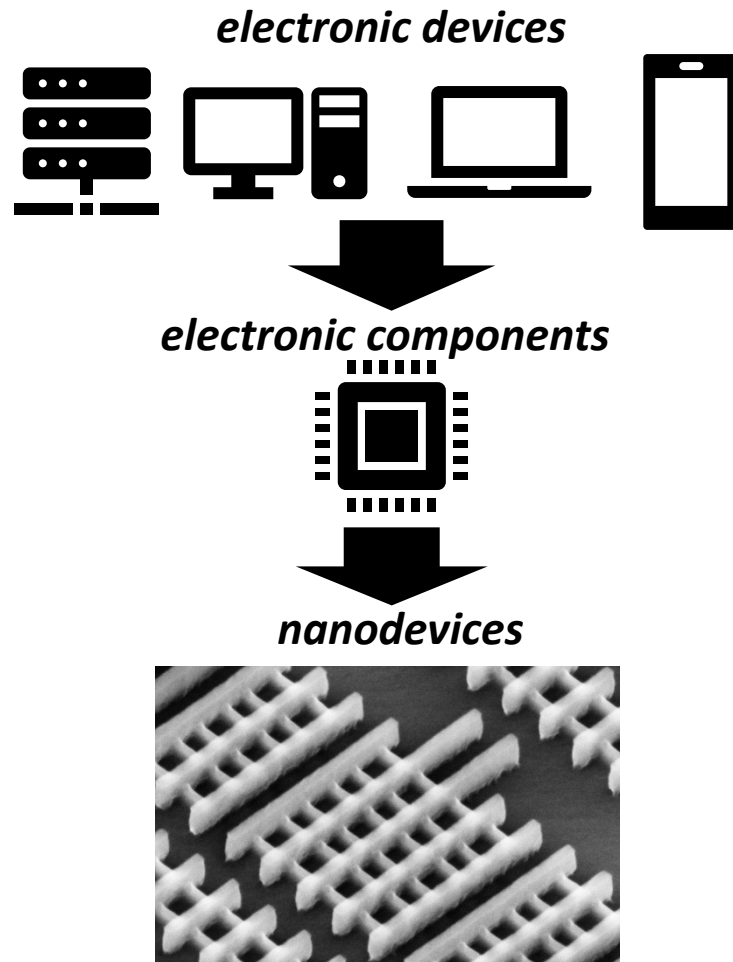A. N. ZIOGAS, TAL BEN-NUN, G. FERNANDÉZ, T. SCHNEIDER, M. LUISIER, T. HOEFLER
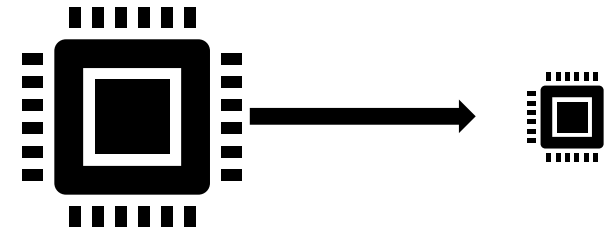
# A Data-Centric Approach to Quantum Transport Simulations

# Motivation
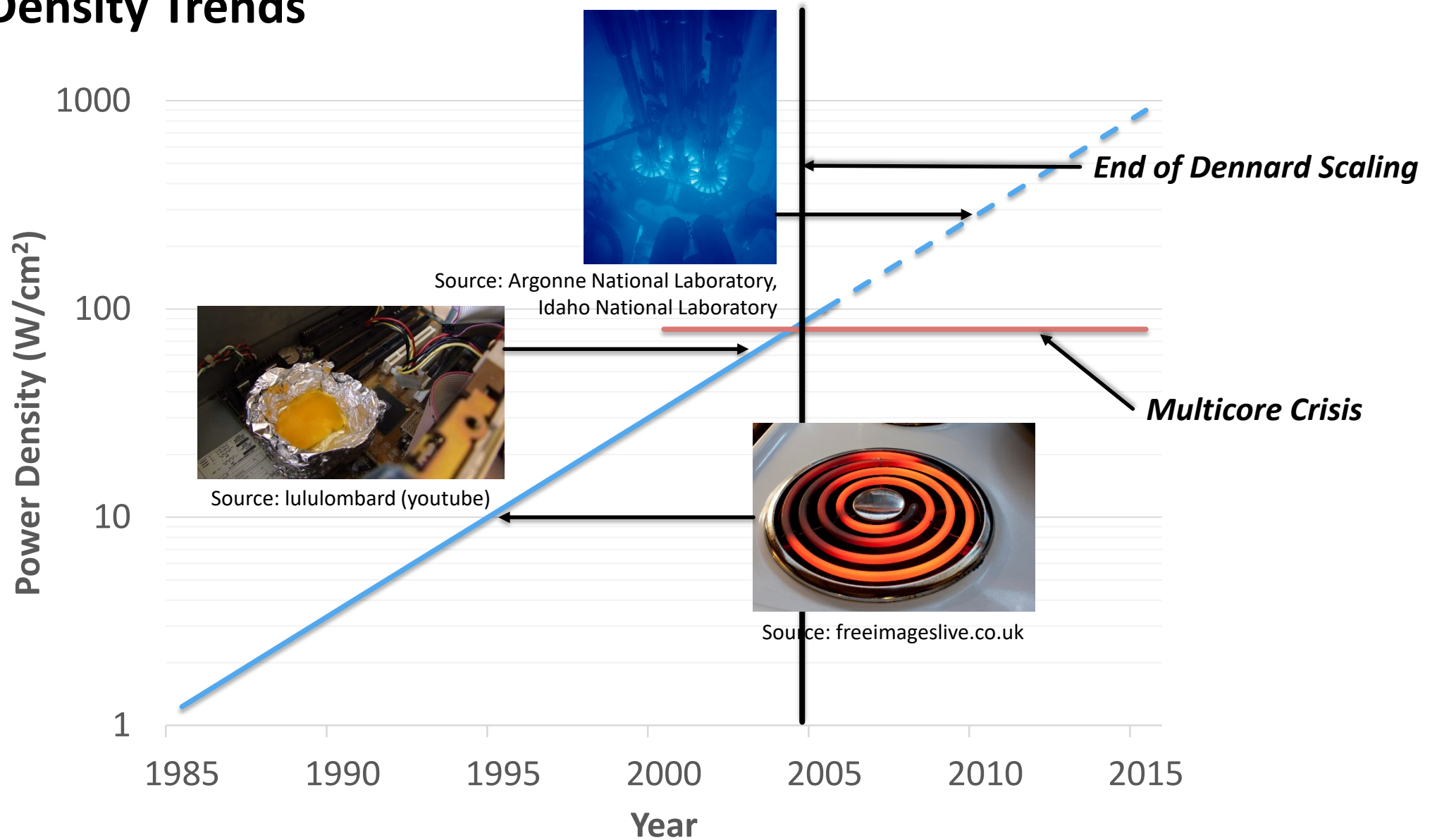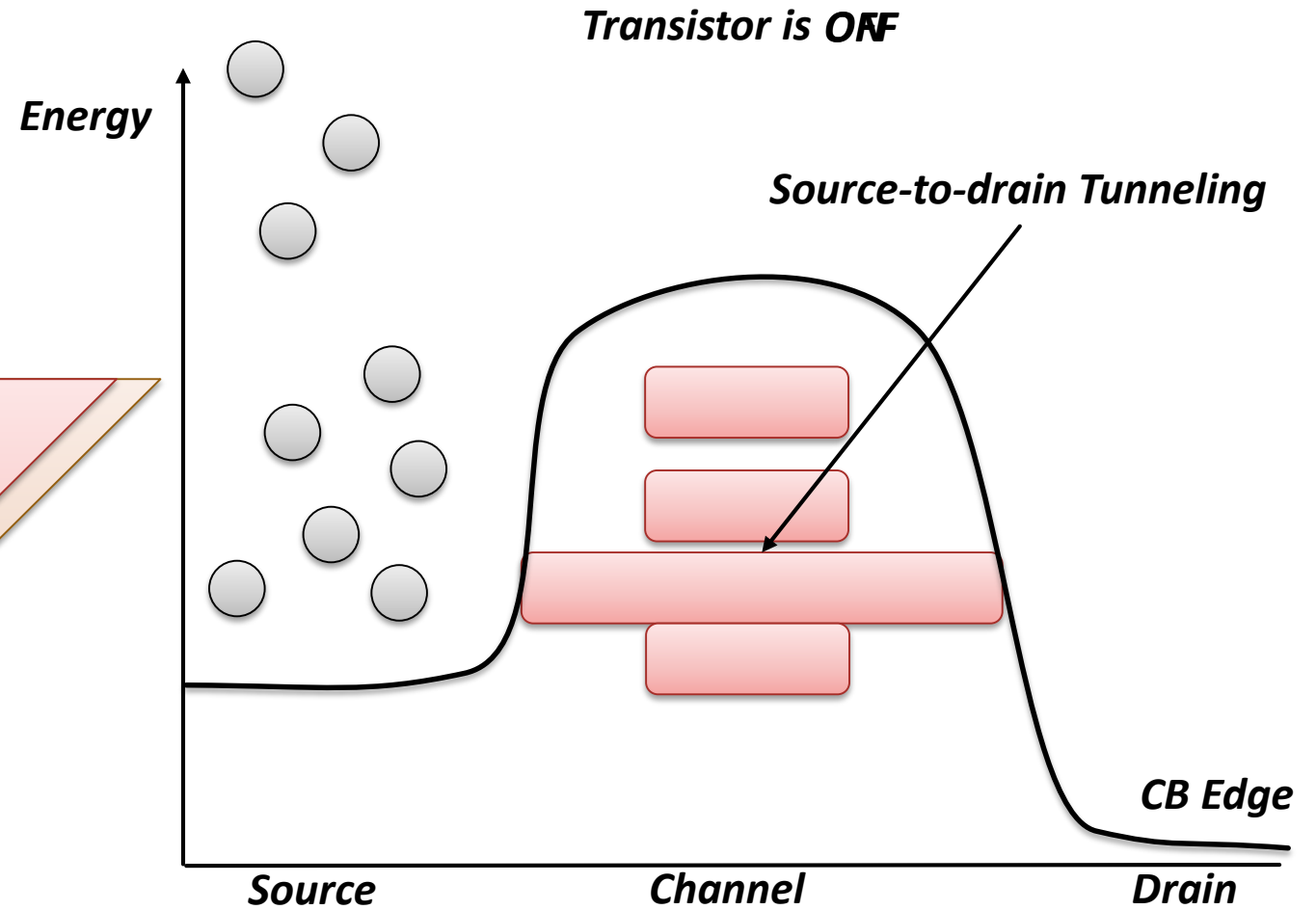
*electronic devices*

*electronic components*

*nanodevices*

Source: Intel

# CPU Power Density Trends



**End of Dennard Scaling**

Source: Argonne National Laboratory,
Idaho National Laboratory

Source: lululombard (youtube)

**Multicore Crisis**

Source: freeimageslive.co.uk
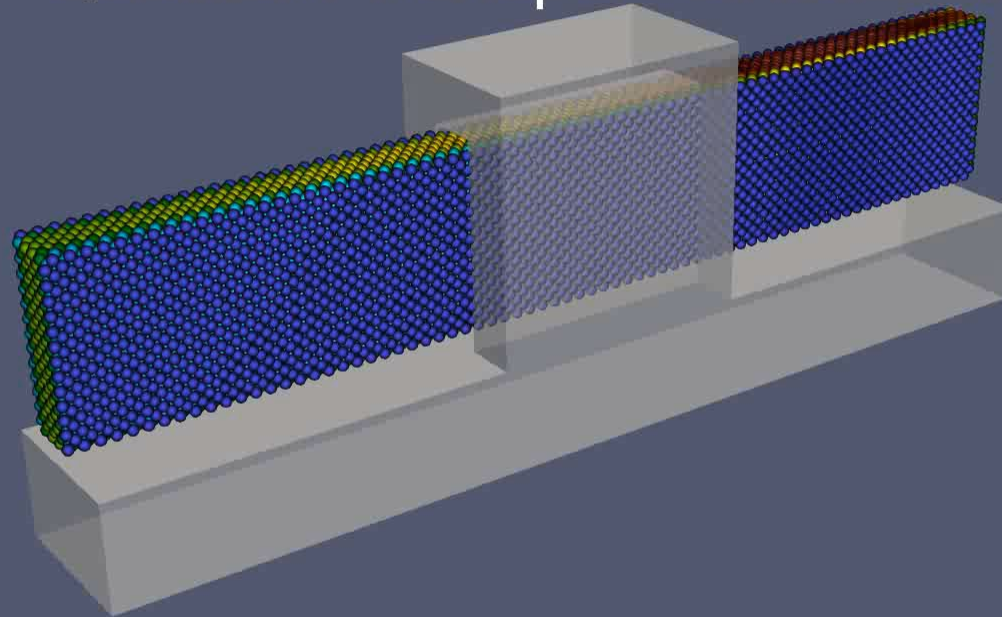
Power Density (W/cm²)
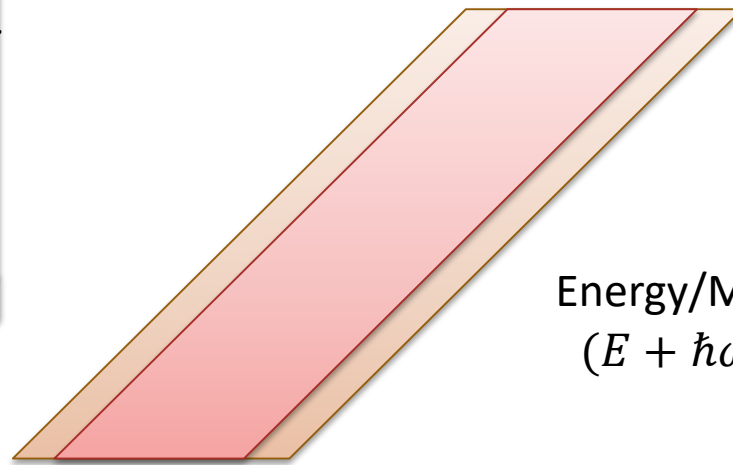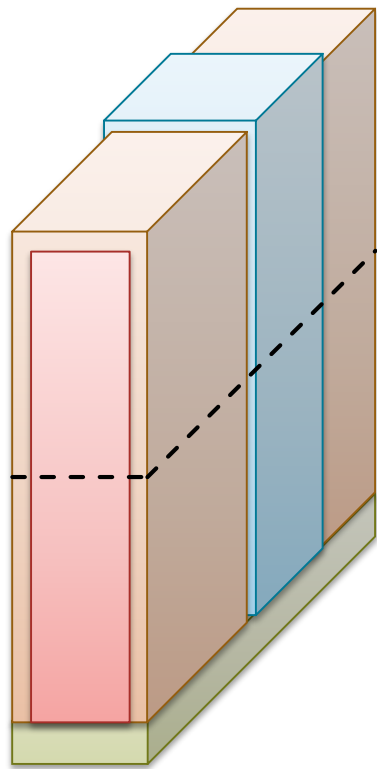
Year

# Quantum Mechanical Phenomena

# Characteristics of Nanotransistors



Extreme-Scale Ab initio Dissipative Quantum Transport Simulations

Alexandros Nikolaos Ziogas, Tal Ben-Nun
Guillermo Indalecio Fernández, Timo Schneider
Mathieu Luisier and Torsten Hoefler

# Quantum Transport Simulation

**Ballistic Transport**

Energy/Momentum

$(E, k)$

Channel

**Incoherent Transport**

$(\omega_1, q_1)$

Energy/Momentum

$(E + \hbar\omega_1, k - q_1)$

$(E, k)$
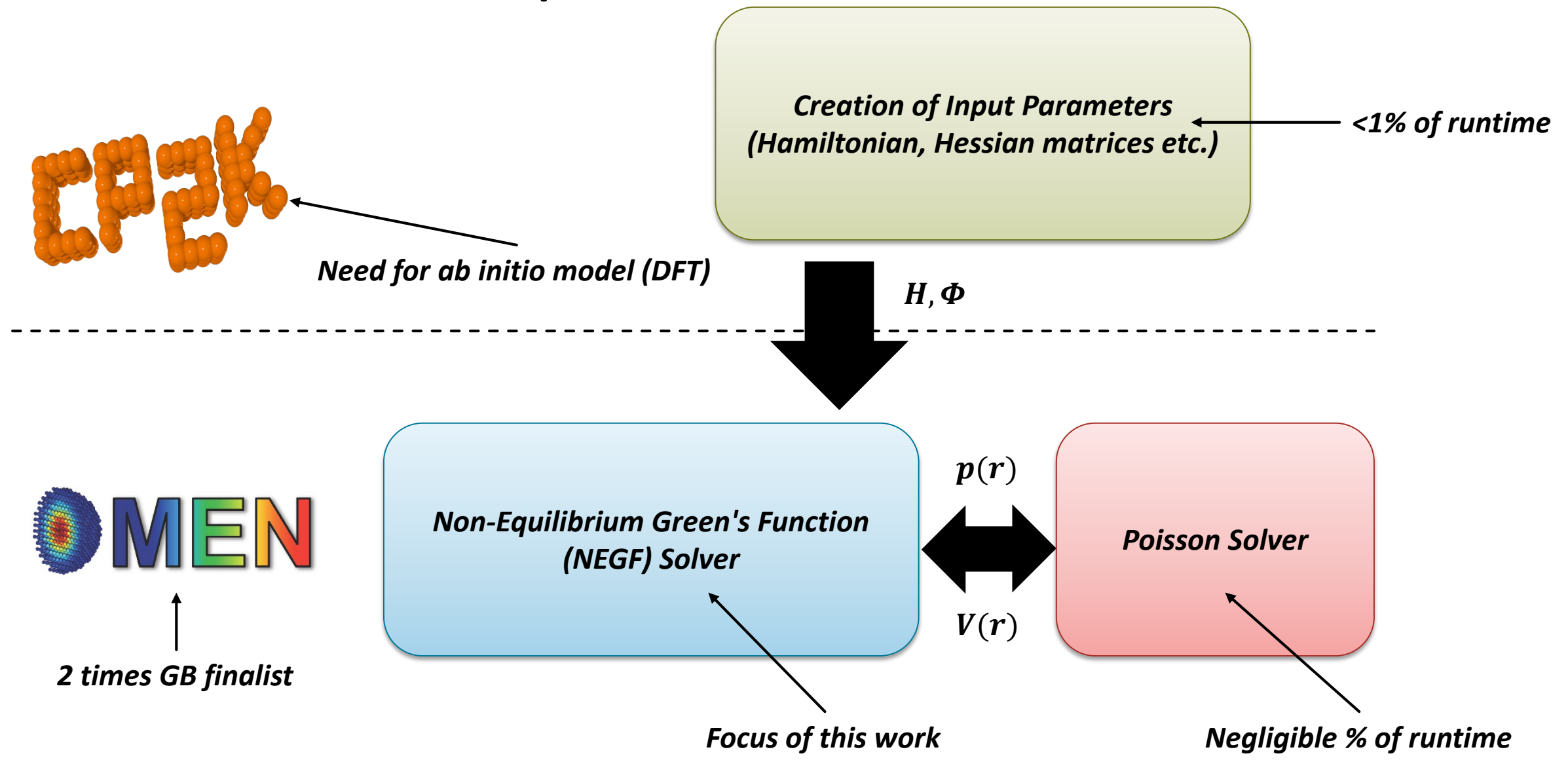
$(E + \hbar(\omega_1 - \omega_2), k - q_1 - q_2)$
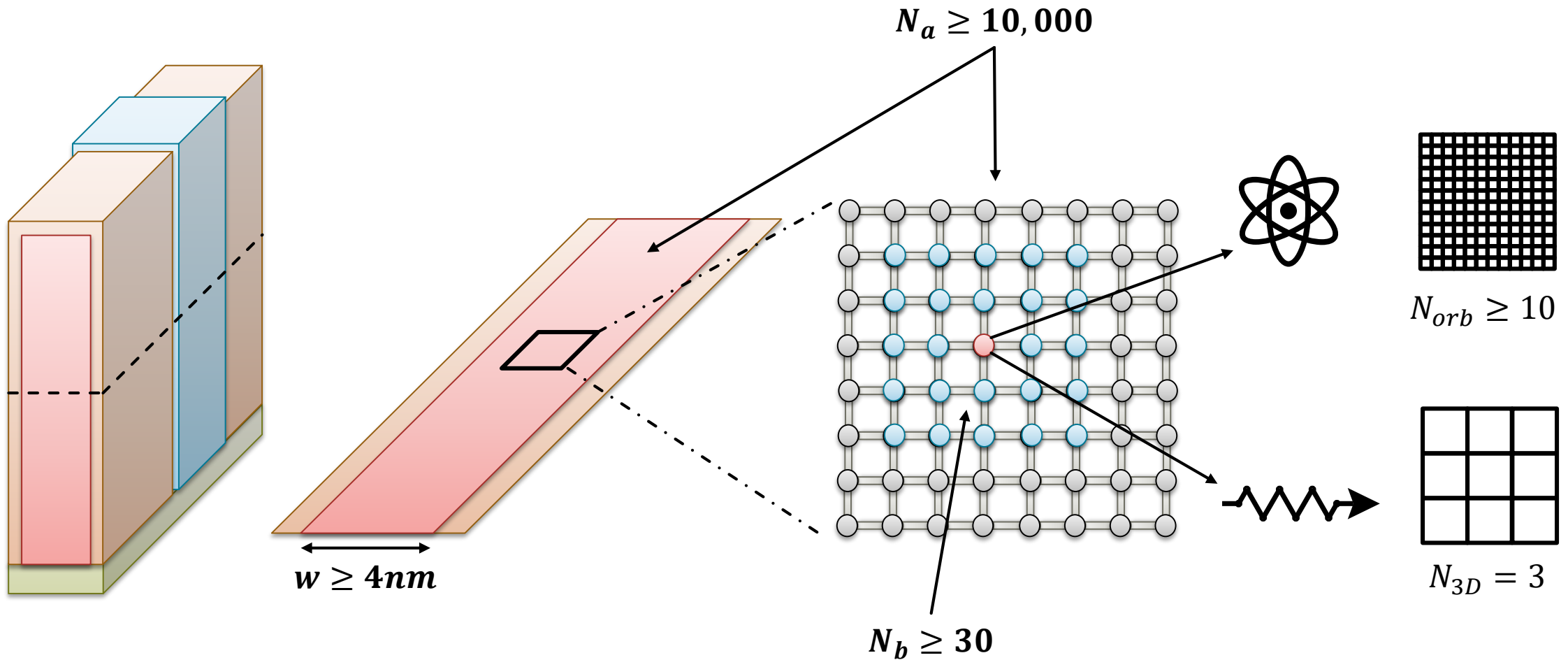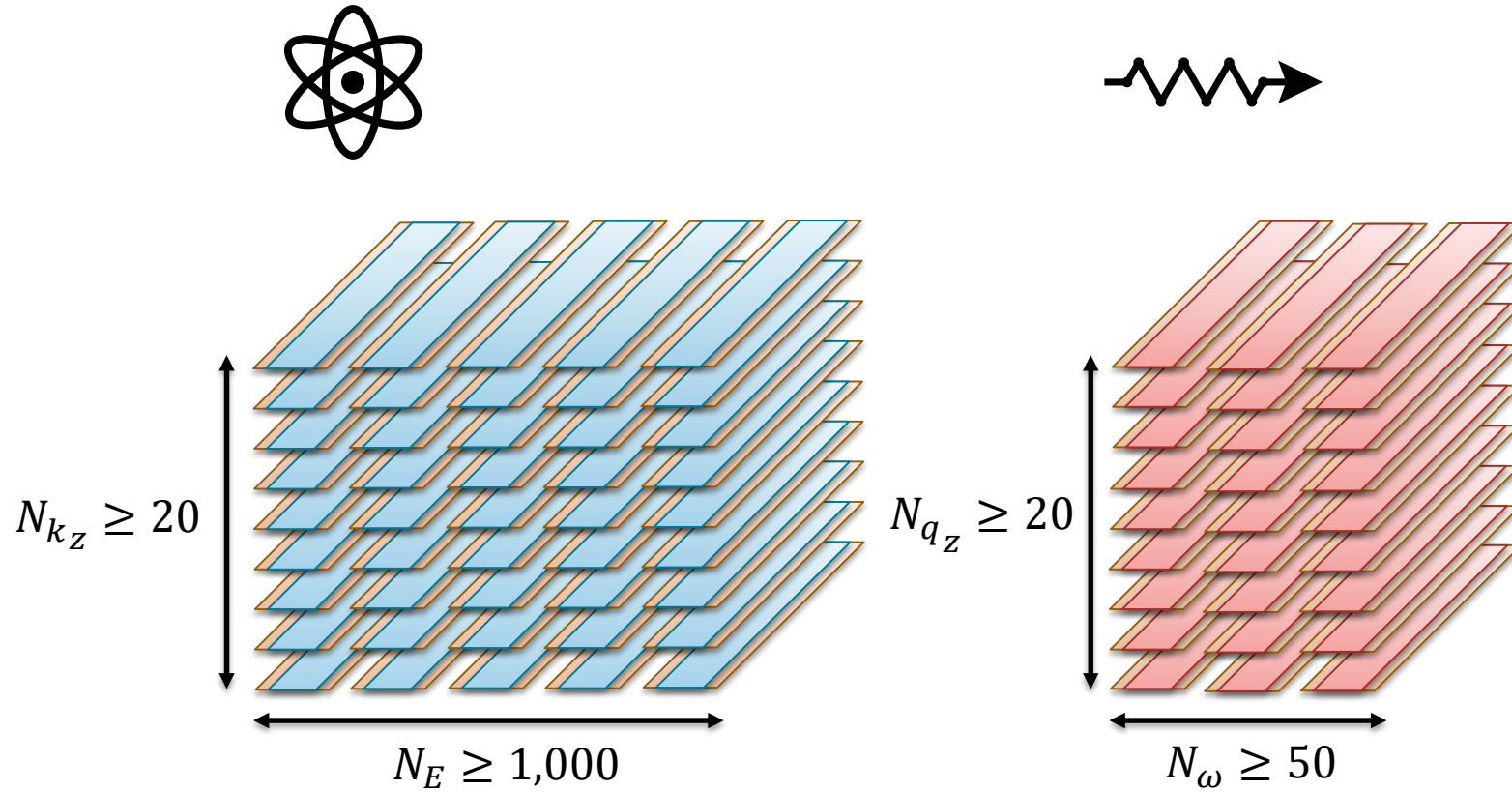
$(\omega_2, q_2)$

**Phonons: Crystal Vibrations**

Channel

6

# Workflow of Quantum Transport Simulation

**Creation of Input Parameters (Hamiltonian, Hessian matrices etc.)**

*<1% of runtime*

*Need for ab initio model (DFT)*

$H, \Phi$

**Non-Equilibrium Green's Function (NEGF) Solver**

$p(r)$

**Poisson Solver**

$V(r)$

*2 times GB finalist*

*Focus of this work*

*Negligible % of runtime*

# OMEN Model



$$N_a \geq 10,000$$

$$w \geq 4nm$$

$$N_b \geq 30$$

$$N_{orb} \geq 10$$

$$N_{3D} = 3$$

# OMEN Model



$$N_{k_z} \geq 20$$

$$N_E \geq 1,000$$

$$N_{q_z} \geq 20$$

$$N_\omega \geq 50$$

# NEGF Mathematical Formulation

**NEGF** $\text{SSE } \Sigma[G(E \pm \hbar\omega, k_z - q_z)\, D(\omega, q_z)](E, k_z)$

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

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

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

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

GF

$\text{SSE } \Pi[G(E, k_z)\, G(E + \hbar\omega, k_z + q_z)](\omega, q_z)$
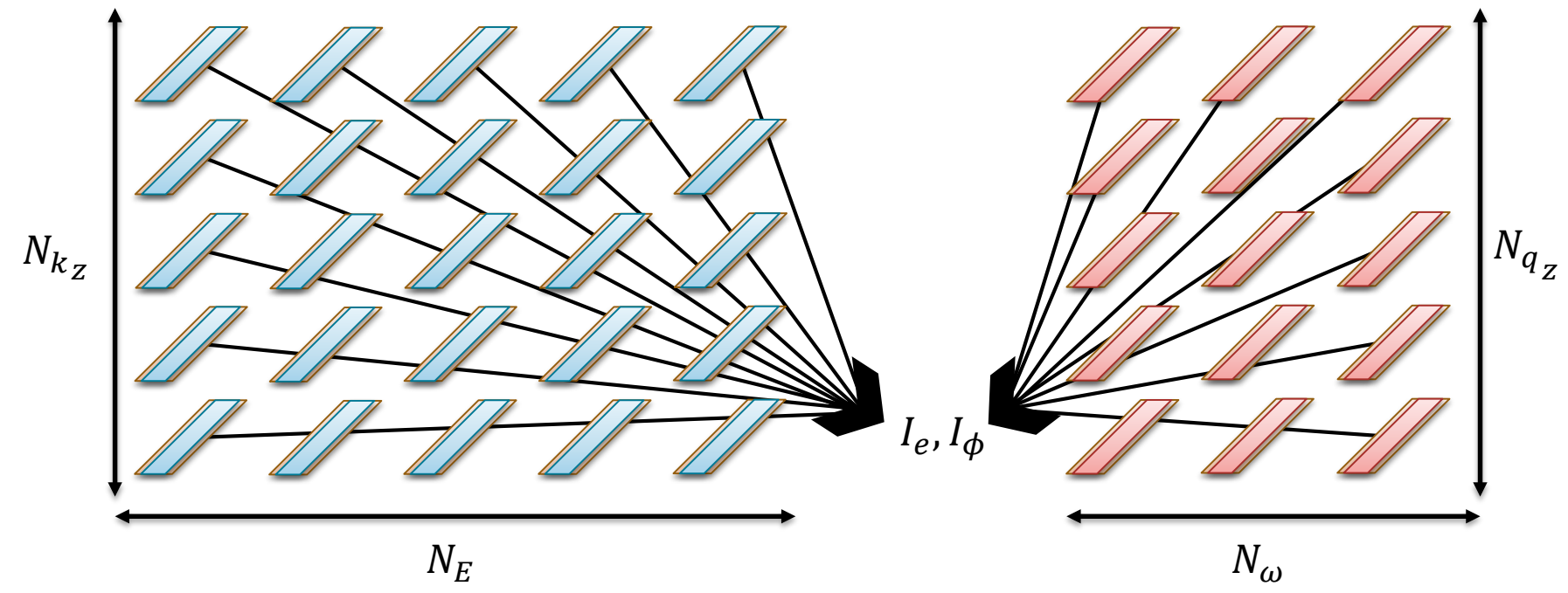
# Assignment to Compute Resources

# Green's Functions Phase



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

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

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

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

GF

*Embarrassingly parallel computation of $G/D$ + Reduction for $I$*
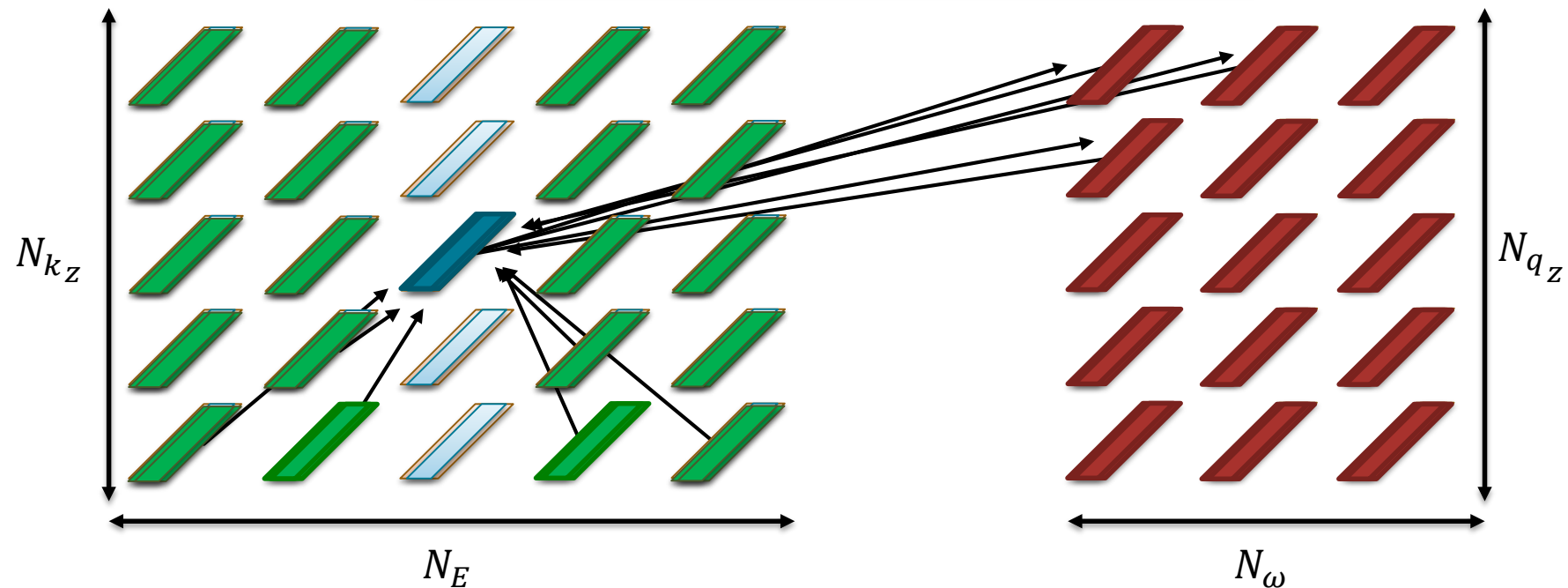
$N_{k_z}$

$N_{q_z}$

$I_e, I_\phi$

$N_E$

$N_\omega$

# Scattering Self-Energies Phase

SSE $\Sigma[G(E \pm \hbar\omega, k_z - q_z) D(\omega, q_z)](E, k_z)$
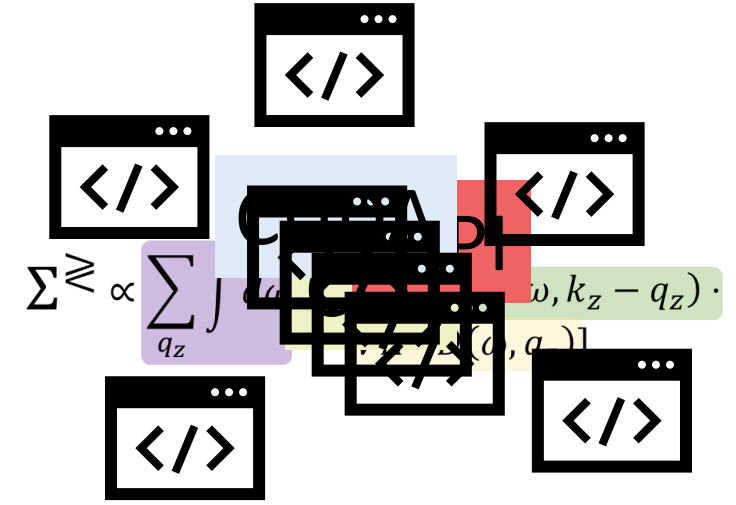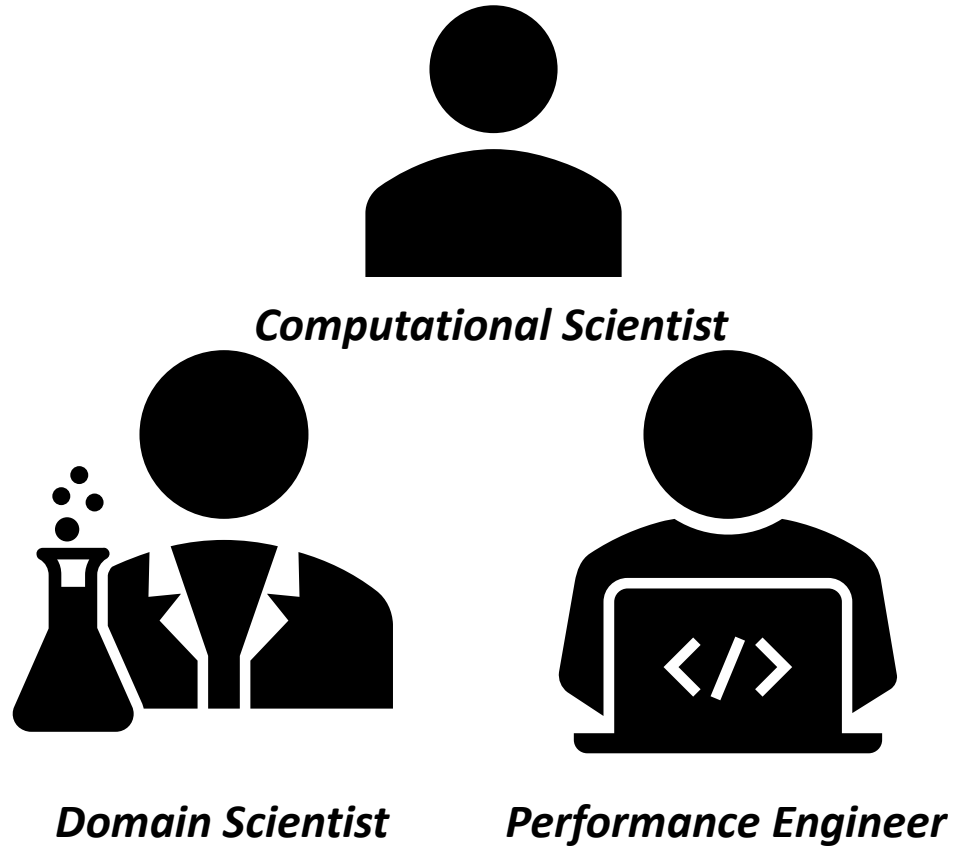
**Stencil-like computation for $\Sigma/\Pi$**
$2N_{q_z}N_\omega$

SSE $\Pi[G(E, k_z) G(E + \hbar\omega, k_z + q_z)](\omega, q_z)$
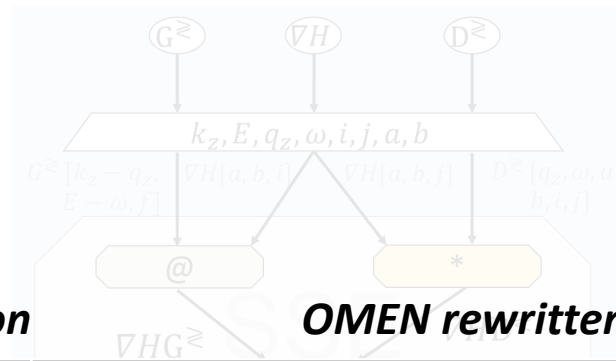
# Motivation for a Data-Centric Framework



*Computational Scientist*

*Domain Scientist*     *Performance Engineer*

$$\Sigma^{\gtrless} \propto \sum_{q_z} \int \dots (\omega, k_z - q_z) \cdot \dots (\omega, q_z)]$$

# Data-Centric Framework



*Original Application*

*OMEN rewritten with the DaCe framework*

| OMEN | | DaCe OMEN |
|---|---|---|
| 15,798 C++ Lines | 3,155 Python Lines | 2,015 SDFG Nodes |

*Data-Centric Intermediate Representation (SDFG)*

*Problem Formulation*

```
@dace.program
def sse_sigma(neigh_idx: dace.int32[NA, NB],
              dH: dace.complex128[NA, NB, N3D, Norb, Norb],
              G: dace.complex128[Nkz, NE, NA, Norb, Norb],
              D: dace.complex128[Nqz, Nw, NA, NB, N3D, N3D],
              Sigma: dace.complex128[Nkz, NE, NA, Norb, Norb]):

    # Declaration of Map scope
    for k, E, q, w, i, j, a, b in dace.map[0:Nkz, 0:NE,
                                           0:Nqz, 0:Nw,
                                           0:N3D, 0:N3D,
                                           0:NA, 0:NB]:
        f = neigh_idx[a, b]
        dHG = G[k-q, E-w, f] @ dH[a, b, i]
        dHD = dH[a, b, j] * D[q, w, a, b, i, j]
        Sigma[k, E, a] += dHG @ dHD
```

**High-Level Program**

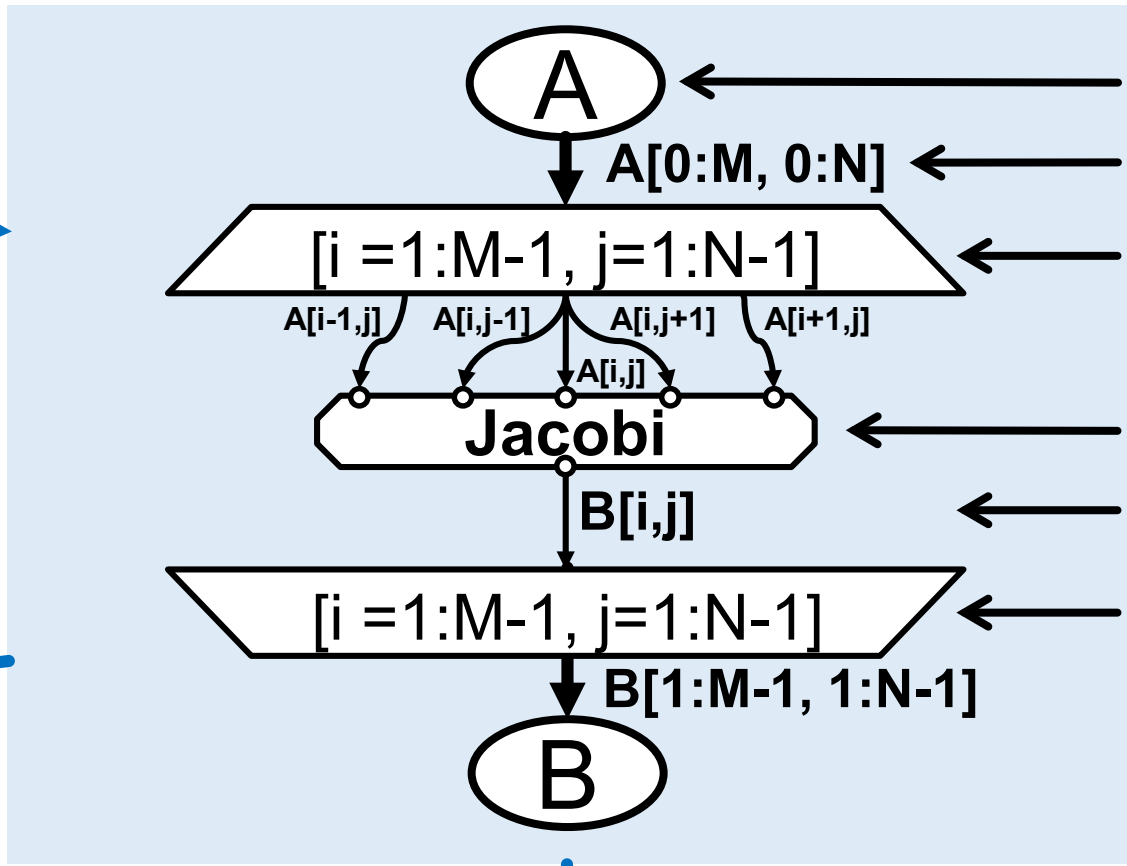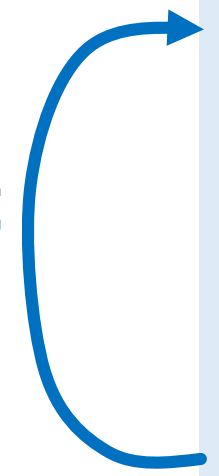*Code Generation*

*Graph Transformations*
*(API, Interactive)*

# Basic Elements of the Data-Centric IR

*Data*

*Fine-grained stateless computations*

*Dataflow*

*Control-flow*

*Abstraction of independently parallel computations*

$t < T;$
$t=t+1$



Access Node

A[0:M, 0:N] — Memlet

[i =1:M-1, j=1:N-1] — Map Entry

A[i-1,j]    A[i,j-1]    A[i,j+1]    A[i+1,j]
A[i,j]

Jacobi — Tasklet

B[i,j] — State

[i =1:M-1, j=1:N-1] — Map Exit

B[1:M-1, 1:N-1]

B

$t \geq T$ — Control flow

Ben-Nun et al. "Stateful Dataflow Multigraphs: A Data-Centric Model for Performance Portability on Heterogeneous Architectures", SC19.

# Data-Centric Representation of OMEN

# Caching Computations



**independent of i**

**Data-centric Transformation**

3GB per $(k_z, E)$

1GB per $(k_z, E)$

**Caching Schemes:**
- **No Cache**
- **Cache BC**
- **Cache BC + specialize H**

**GF**

**GF**

18

# Optimization of Sparse Operations

```python
for n in range(N - 2, -1, -1):
    sig = HF[n] @ gR[n + 1] @ HE[n + 1]
```

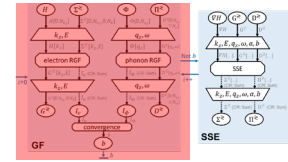*Sparse (either CSR or CSC)*          *Dense*

| Method\Operation | NN | NT | TN | TT |
|---|---|---|---|---|
| cublasZgemm | 58.382 ms | 58.144 ms | 58.666 ms | 58.315 ms |
| cusparseZcsrmm2 | 8.202 ms | 6.140 ms | 52.722 ms | – |
| cusparseZgemmi | 15.198 ms | – | – | – |

| Approach | Time |
|---|---|
| gemm/gemm | 116.881 ms |
| gemmi(csrmm2(TN, HE, gR), HF) | 67.924 ms |
| csrmm2(NT, HE, csrmm2(NT, HF, gR)) | 11.994 ms |

*HF and HE are CSR*

*HF is CSR, HE is CSC*

# Extracting Parallelism



```
gpR[..1] = np.linalg.inv(tmp - sigRn)
```

```cpp
auto __a = dace::ArrayViewIn<dace::complex128, 2, 1> (gpu_tmpL, bsize, 1);
auto *a = __a.ptr<1>();
auto __b = dace::ArrayViewIn<dace::complex128, 2, 1> (gpu_hergR, bsize, 1);
auto *b = __b.ptr<1>();

auto __c = dace::ArrayViewOut<dace::complex128, 2, 1> (gpu_tmpL_R, bsize, 1);
auto *c = __c.ptr<1>();

int __dace_current_stream_id = 2;
cudaStream_t __dace_current_stream = dace::cuda::__streams[__dace_current_stream_id];

cublasSetStream(handle, __dace_current_stream);
cublasStatus_t status = cublasZgemm(
    handle,
    CUBLAS_OP_N, CUBLAS_OP_N,
    bsize, bsize, bsize,
    const_pone,
    (cuDoubleComplex*)b, bsize,
    (cuDoubleComplex*)a, bsize,
    const_zero,
    (cuDoubleComplex*)c, bsize
);

cudaEventRecord(dace::cuda::__events[5], dace::cuda::__streams[2]);
cudaStreamWaitEvent(dace::cuda::__streams[0], dace::cuda::__events[5], 0);
```
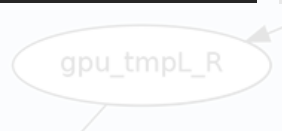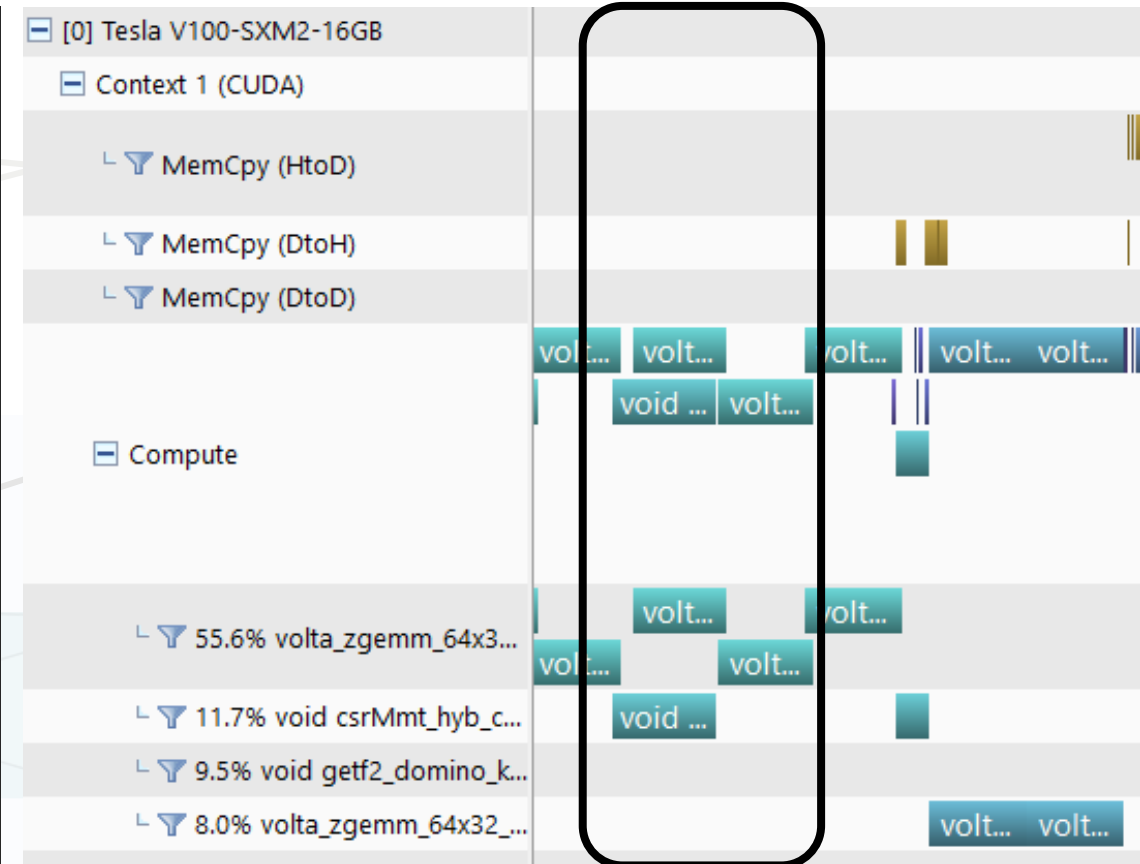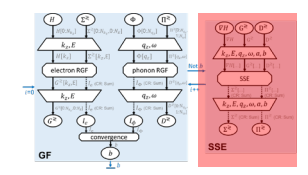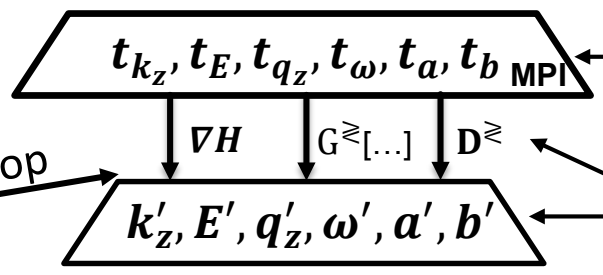
*assignment to streams*

*synchronization code*

gpu_tmpG_R

gpu_tmpL_R

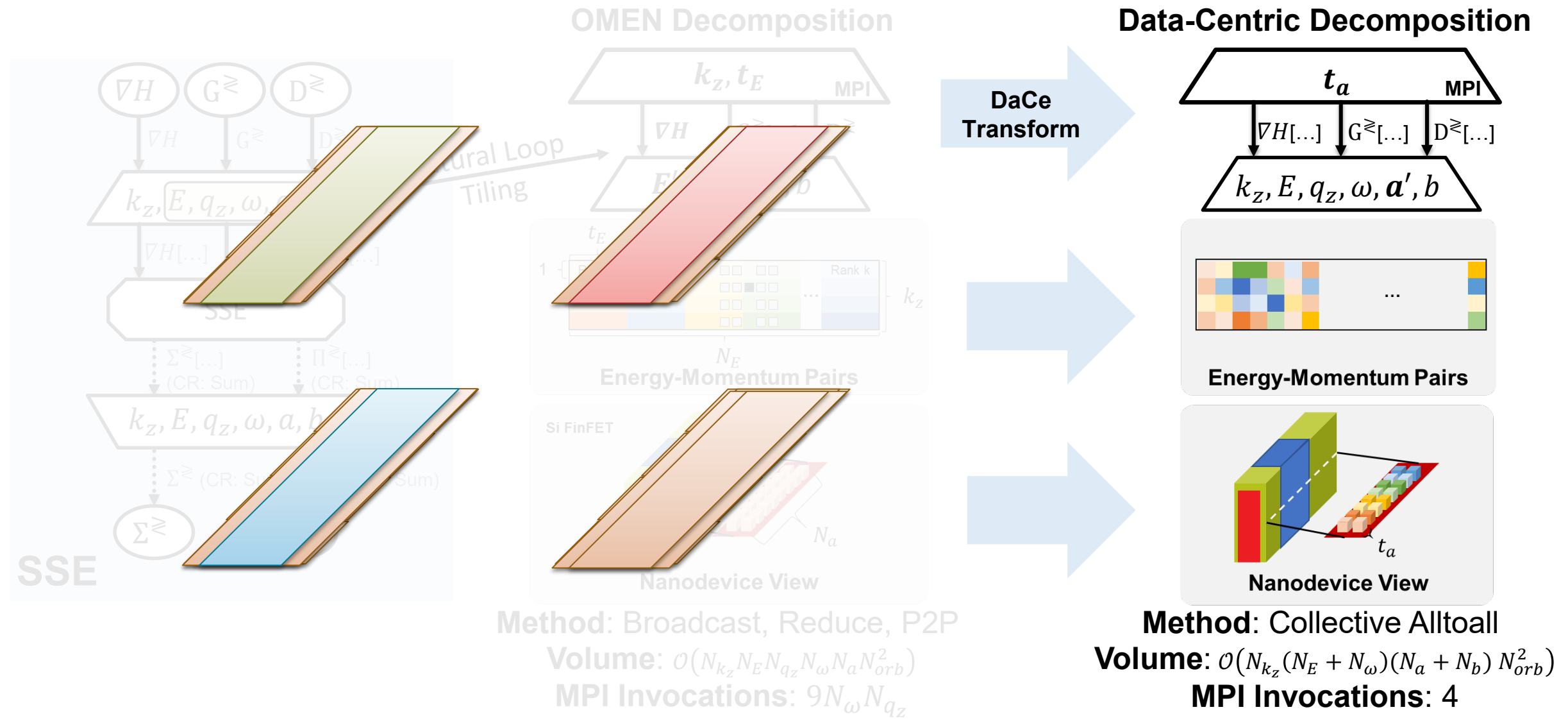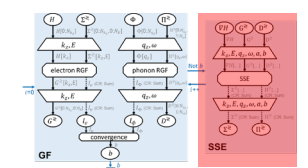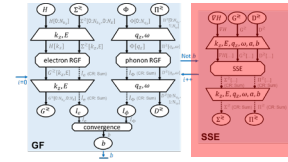# Optimizing Coarse-Grained Data Movement



Iterates over the tiles (nodes)

Iterates over a single tile
(workload of a single node)

Memlets represent the surface of each tile
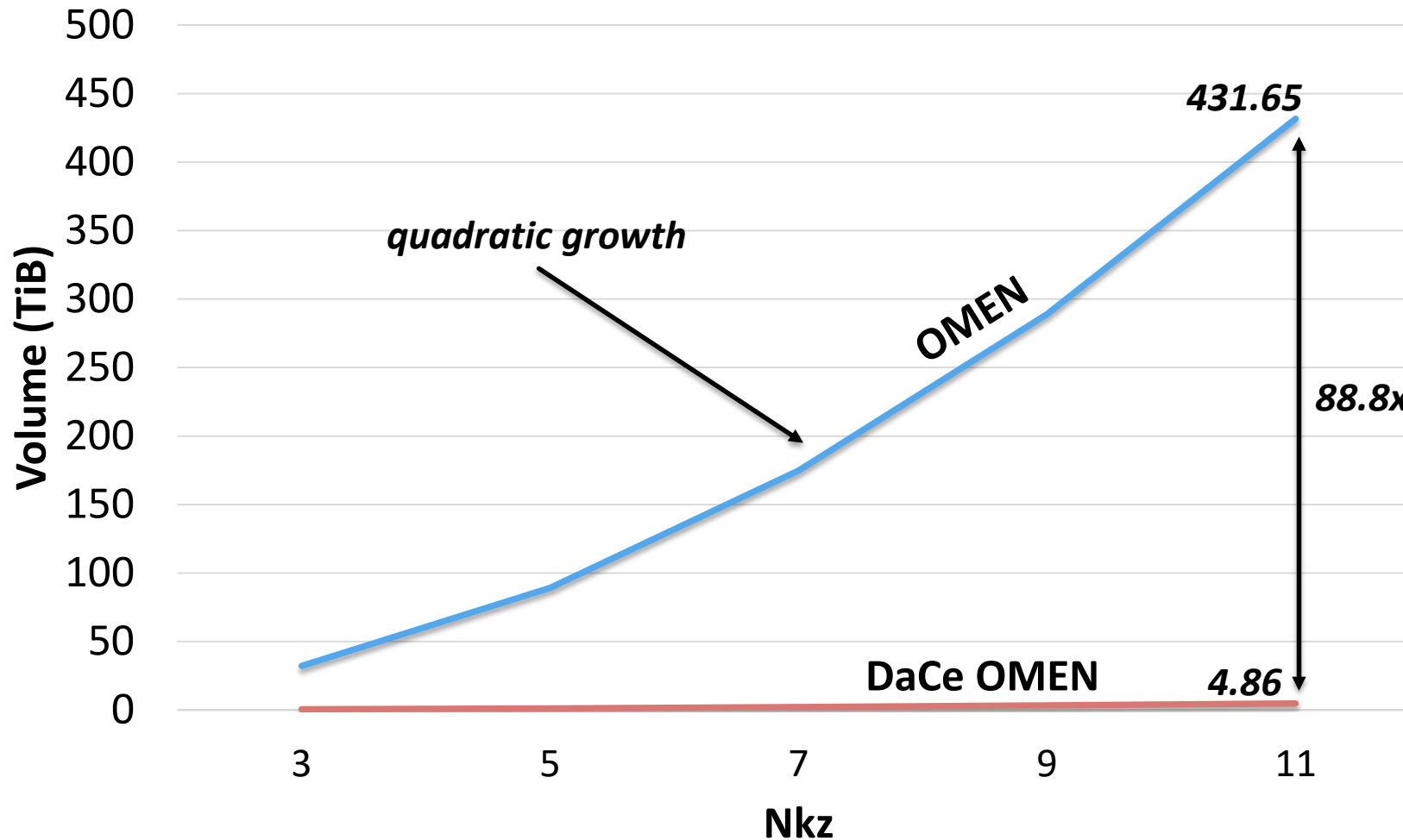(communication volume needed for each node)

Natural Loop Tiling

**SSE**

# Optimizing Coarse-Grained Data Movement



OMEN Decomposition

$k_z, t_E$  MPI

Energy-Momentum Pairs

Si FinFET

Nanodevice View

**Method**: Broadcast, Reduce, P2P
**Volume**: $\mathcal{O}(N_{k_z} N_E N_{q_z} N_\omega N_a N_{orb}^2)$
**MPI Invocations**: $9 N_\omega N_{q_z}$

DaCe Transform

**Data-Centric Decomposition**

$t_a$  MPI

$\nabla H[\dots]$  $G^{\gtrless}[\dots]$  $D^{\gtrless}[\dots]$

$k_z, E, q_z, \omega, \boldsymbol{a'}, b$

**Energy-Momentum Pairs**

$t_a$

**Nanodevice View**

**Method**: Collective Alltoall
**Volume**: $\mathcal{O}(N_{k_z}(N_E + N_\omega)(N_a + N_b) N_{orb}^2)$
**MPI Invocations**: 4

# Optimizing Coarse-Grained Data Movement



Communication Volume

# Optimizing Coarse-Grained Data Movement



Source: Oak Ridge National Laboratory

**Cluster with 100 nodes**

**10k atoms**

**100 atoms per node**

$$t_a, t_E \quad \text{MPI}$$

$$\nabla H[\dots] \quad G^{\gtrless}[\dots] \quad D^{\gtrless}[\dots]$$
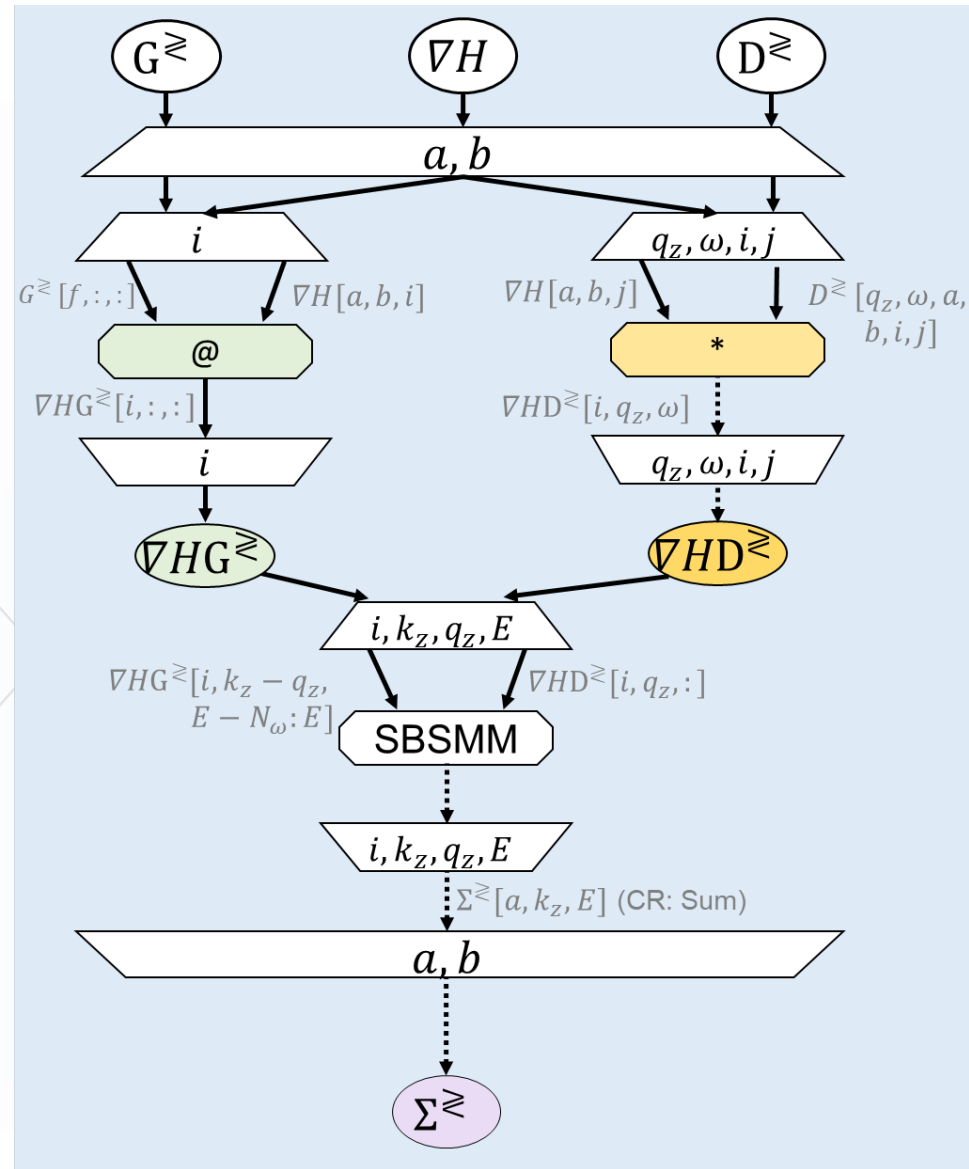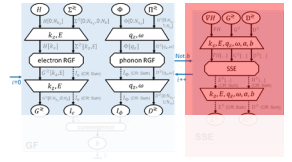
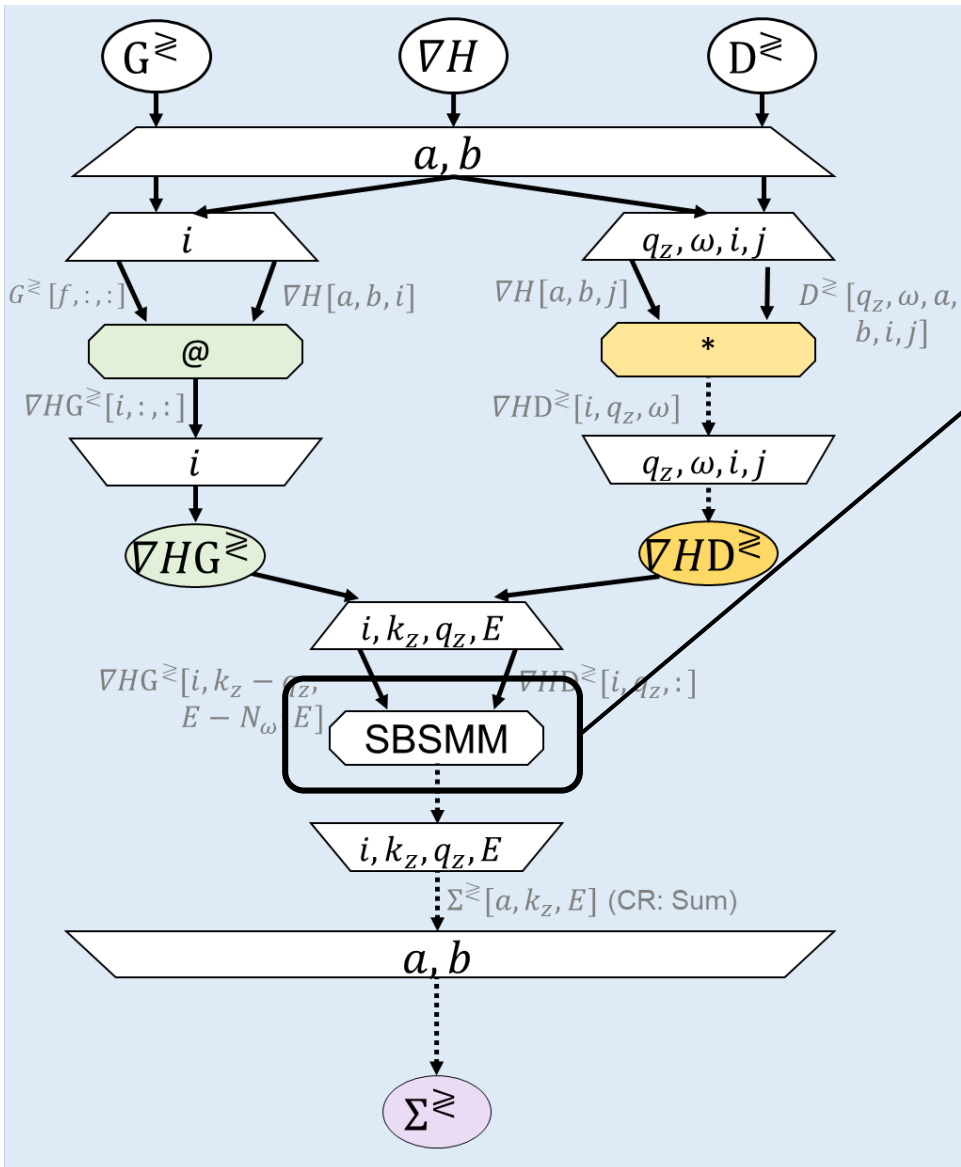$$k_z, E', q_z, \omega, a', b$$
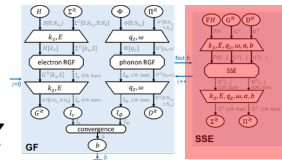
*Need for flexibility*
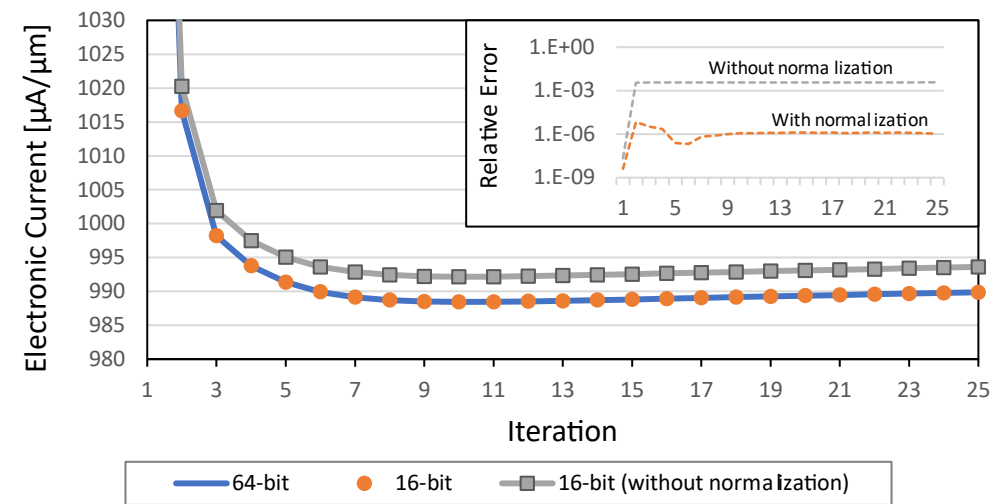
*27k* **X**

Source: NVIDIA
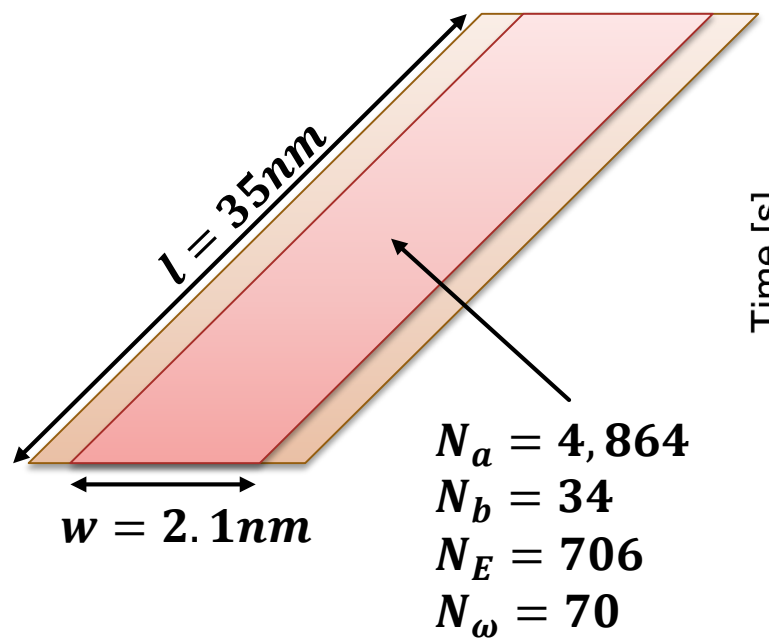
# Optimizing Fine-Grained Data Movement

# Mixed Precision



***cublasGemmStridedBatchedEx***

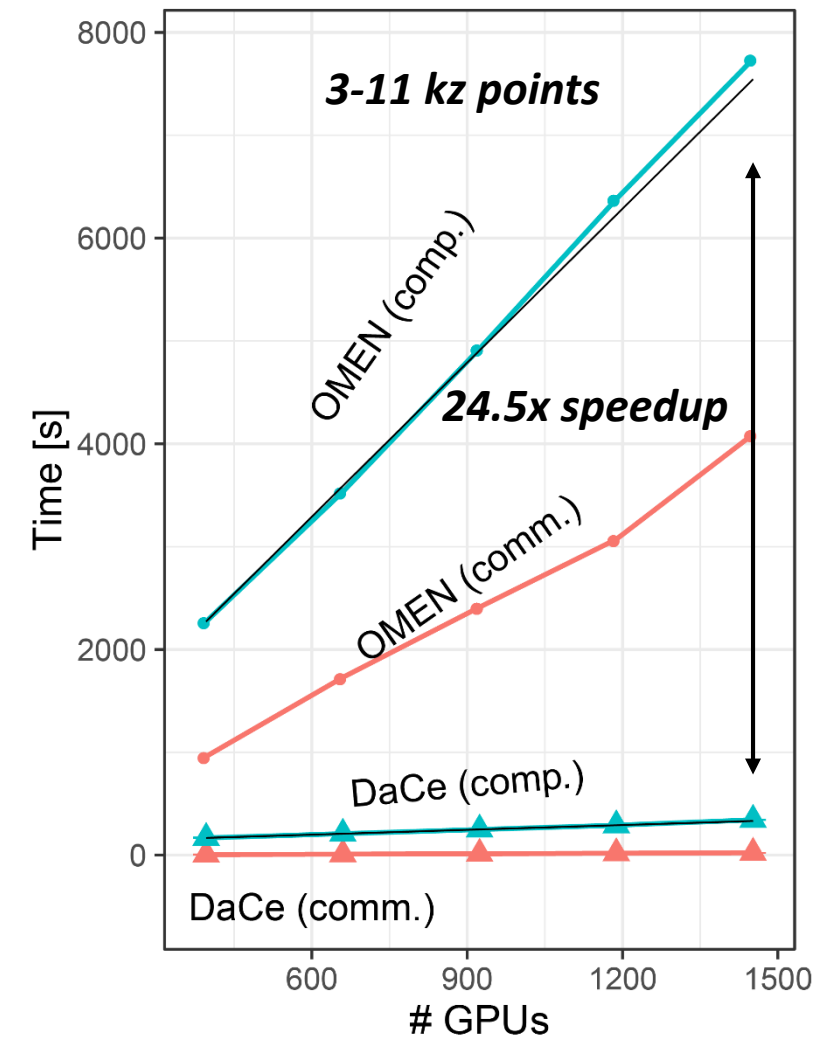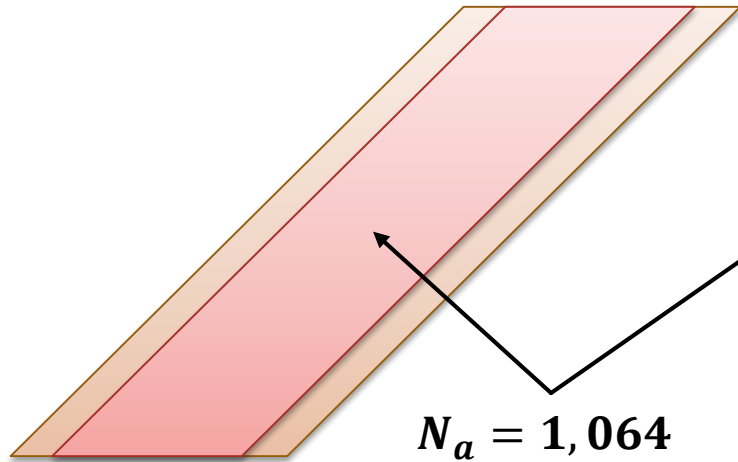| Precision\Implementation | cuBLAS | DaCe |
|---|---|---|
| Double | 4.62 ms | 0.70 ms |
| Mixed – Tensor Cores | N/A | 0.13 ms |

# OMEN vs DaCe OMEN: Performance



Strong Scaling (7 kz points)

Weak Scaling

$l = 35nm$

$w = 2.1nm$

$N_a = 4,864$
$N_b = 34$
$N_E = 706$
$N_\omega = 70$

OMEN (comp.)

OMEN (comm.)

DaCe

(comp.)

(comm.)

**17x speedup**

*3-11 kz points*

OMEN (comp.)

OMEN (comm.)
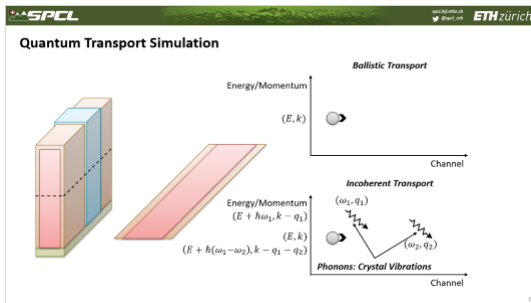
*24.5x speedup*

DaCe (comp.)

DaCe (comm.)

# Beyond OMEN



$N_a = 10,240$
$N_b = 34$
$N_{k_z} = 21$
$N_E = 1220$
$N_\omega = 70$

# Beyond OMEN

**6,840 GPUs**

| Variant | Atom No. | Time [s] | Time/Atom [s] | Speedup |
|---------|---------:|---------:|--------------:|--------:|
| OMEN | 1,064 | 4,695.70 | 4.413 | |
| DaCe OMEN | 10,240 | 333.36 | 0.033 | 140.9x |

$N_a = 1,064$
$N_b = 34$
$N_{k_z} = 21$
$N_E = 1220$
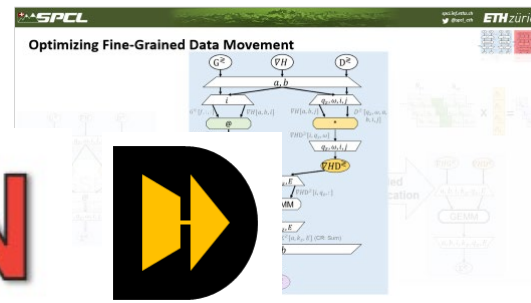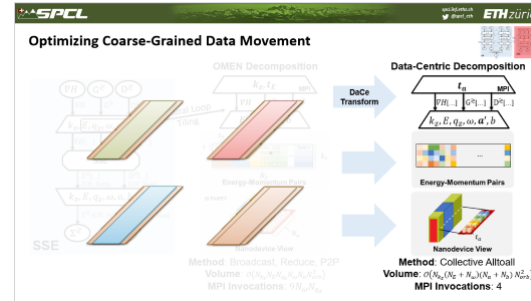$N_\omega = 70$

# Conclusions
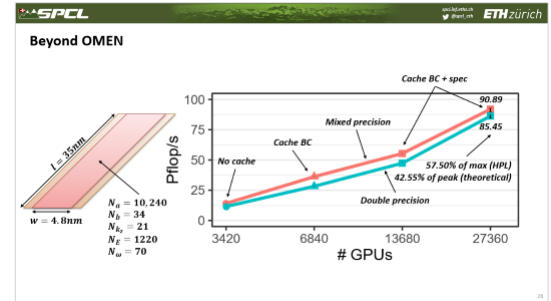
**Heat Dissipation Issue**

**QT Simulation**

**OMEN Application**

**Domain Scientists' View**

**Data-Centric View**

**Optimizing Coarse-Grained and Fine-Grained Dataflow**

**Extracting Parallelism**

**Performance**



github.com/spcl/dace