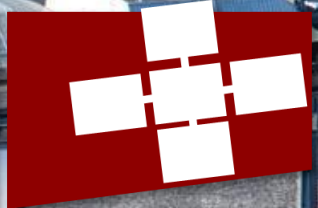


ASA: Application-Specific Architecture

Tiziano De Matteis, Torsten Hoefler

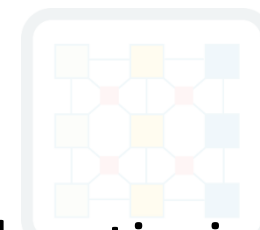


Motivations

Domain Specific Architectures



Spatial Specialized Hardware (SoC)

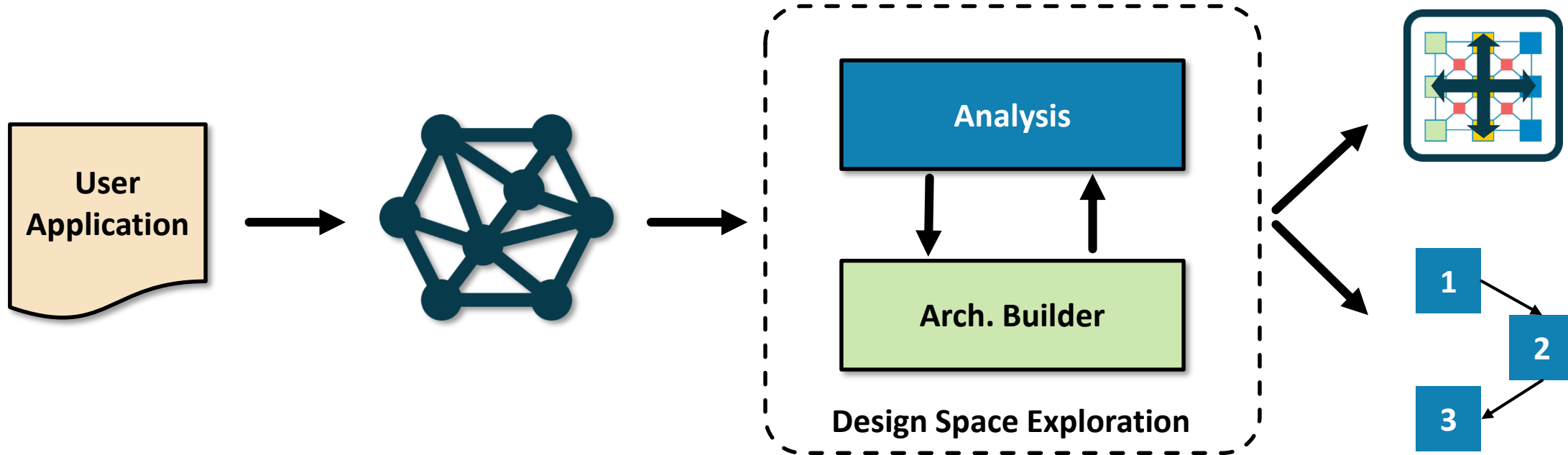


Instead of having an application being adapted and optimized for a given architecture (DSA or SoC), **how we can build a Spatial Specialized Architecture to efficiently support a given application? How we can map the application on such architecture?**

High specialization can be a problem:
algorithms can evolve rapidly making DSA less efficient

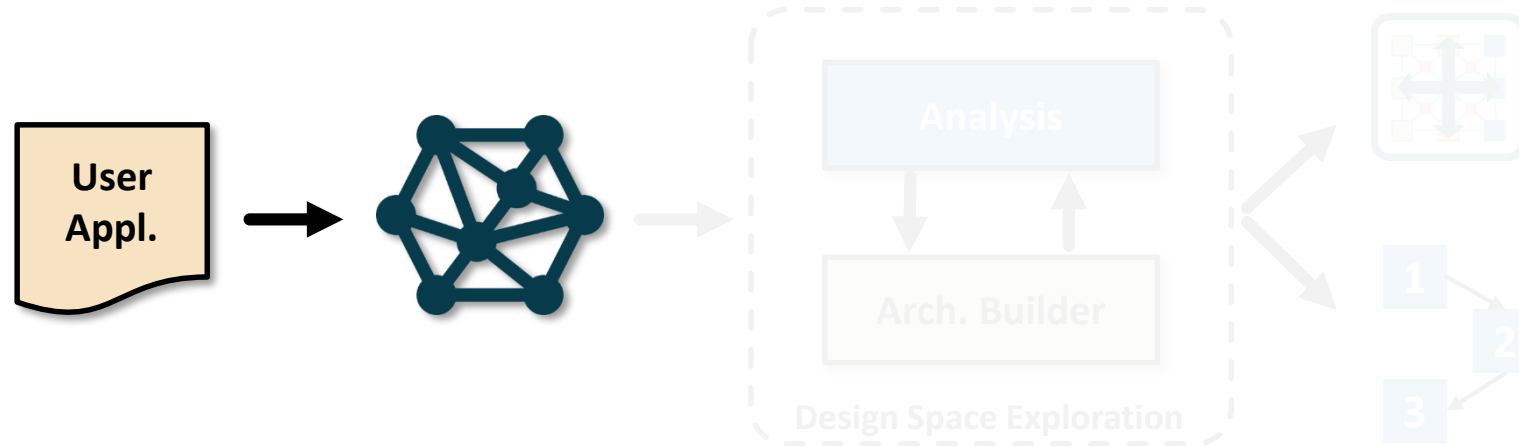
Mapping applications to the architecture is
a complicated task

EFCL Proposal



SDRs and Deep Learning as driving use cases

First steps



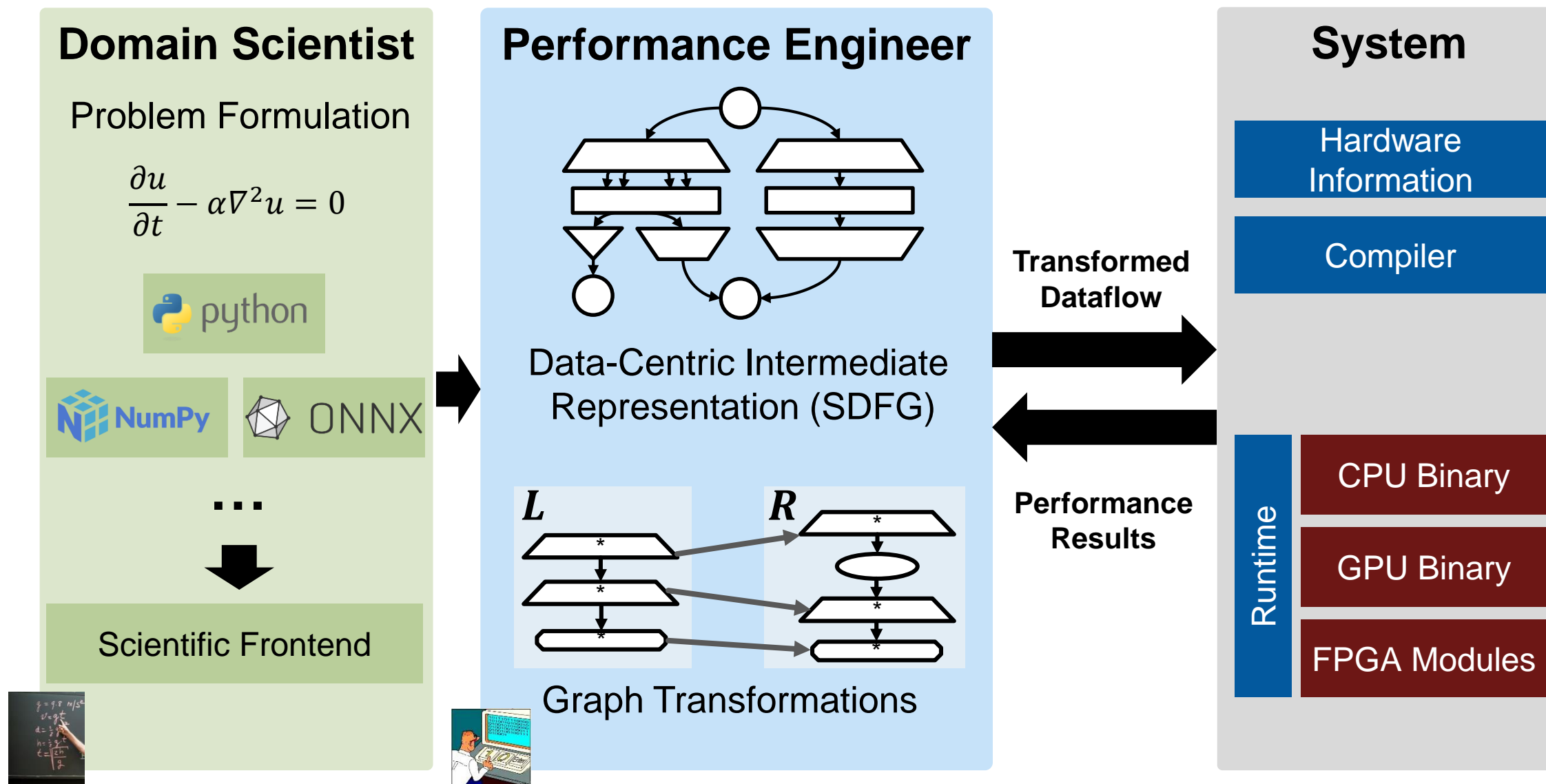
To facilitate further analysis, we want to use data-flow languages to represent the application:

- that naturally describe the **business logic** of the application
- and its **data movements**

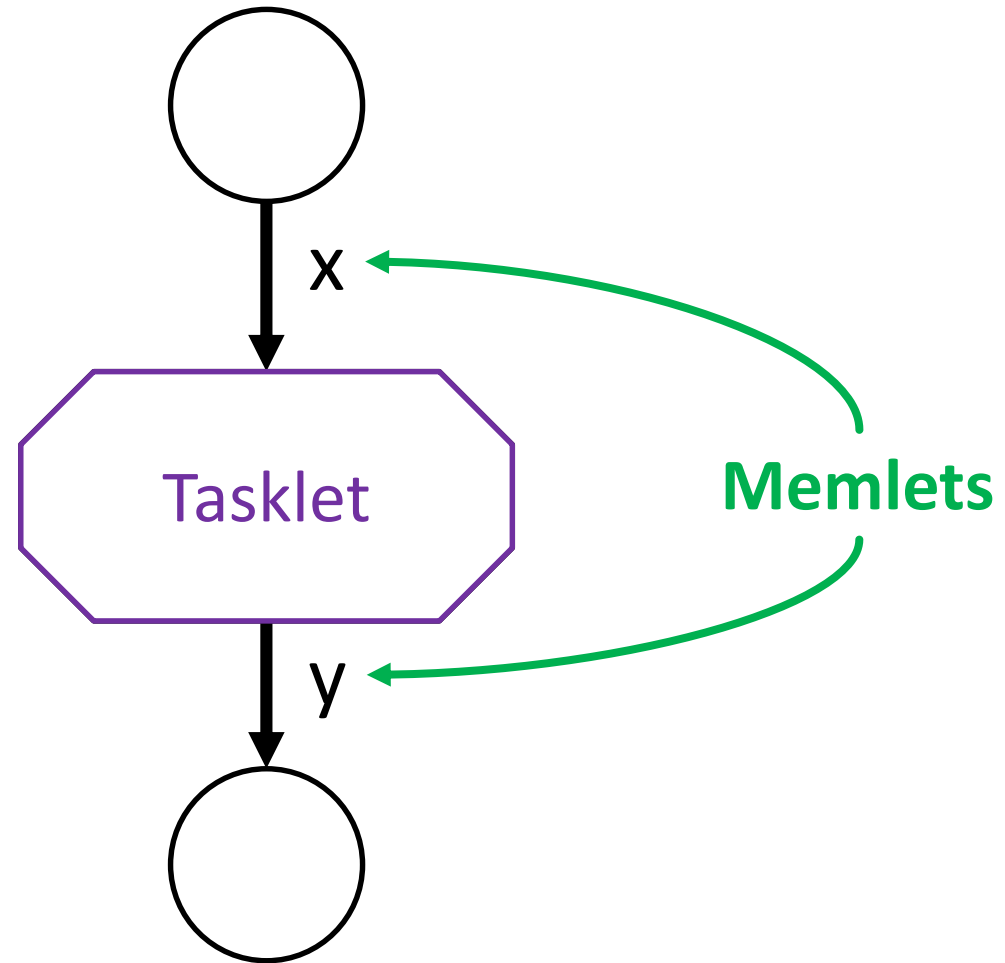
We want to exploit **both** these types of information



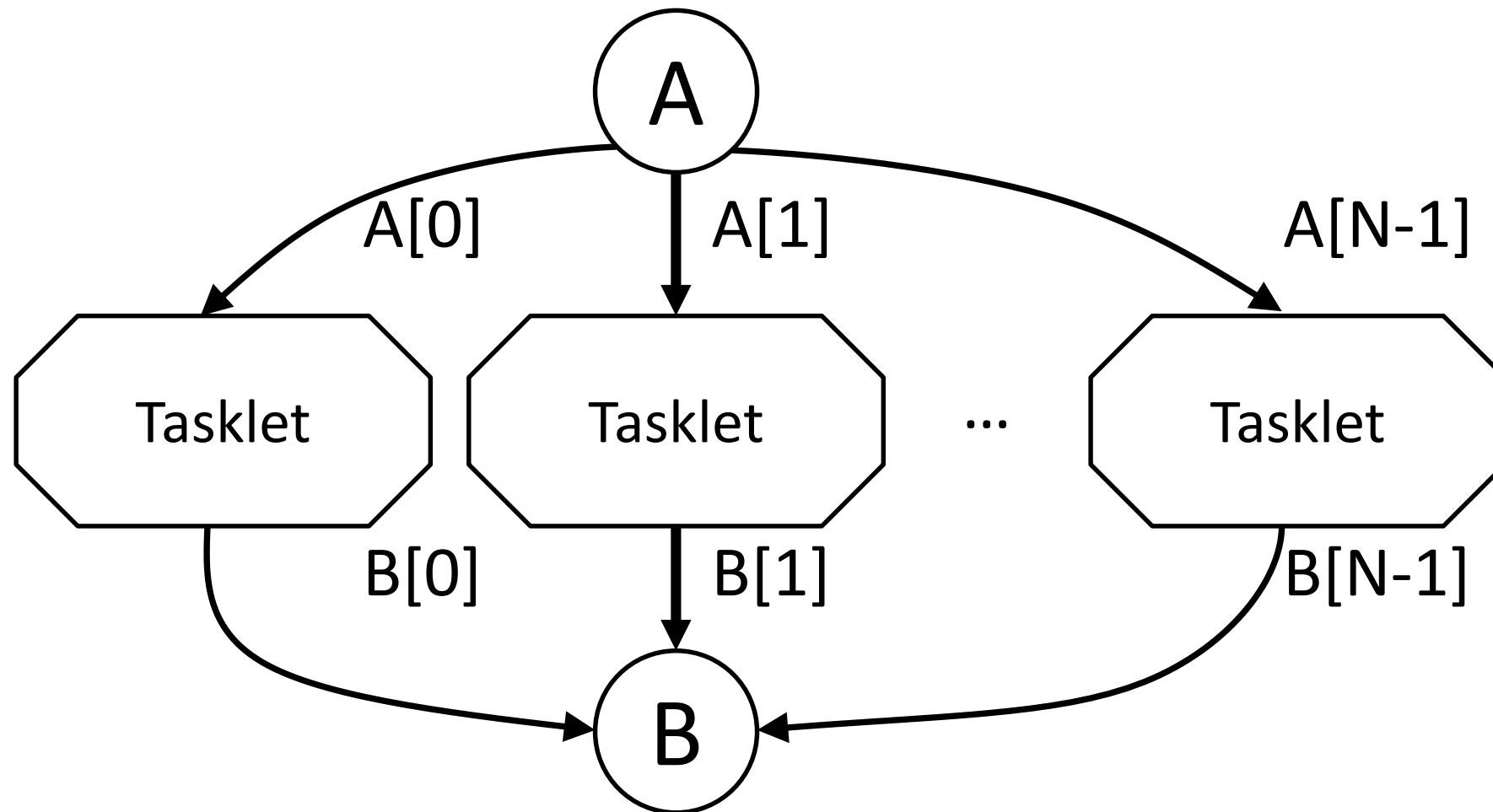
DaCe Overview



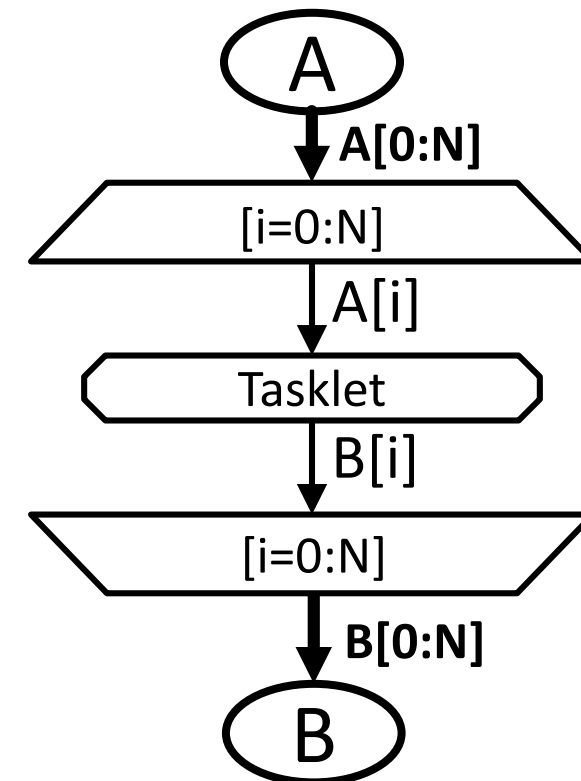
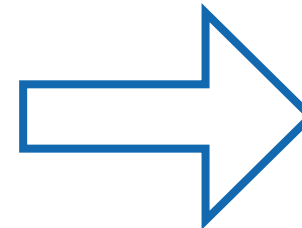
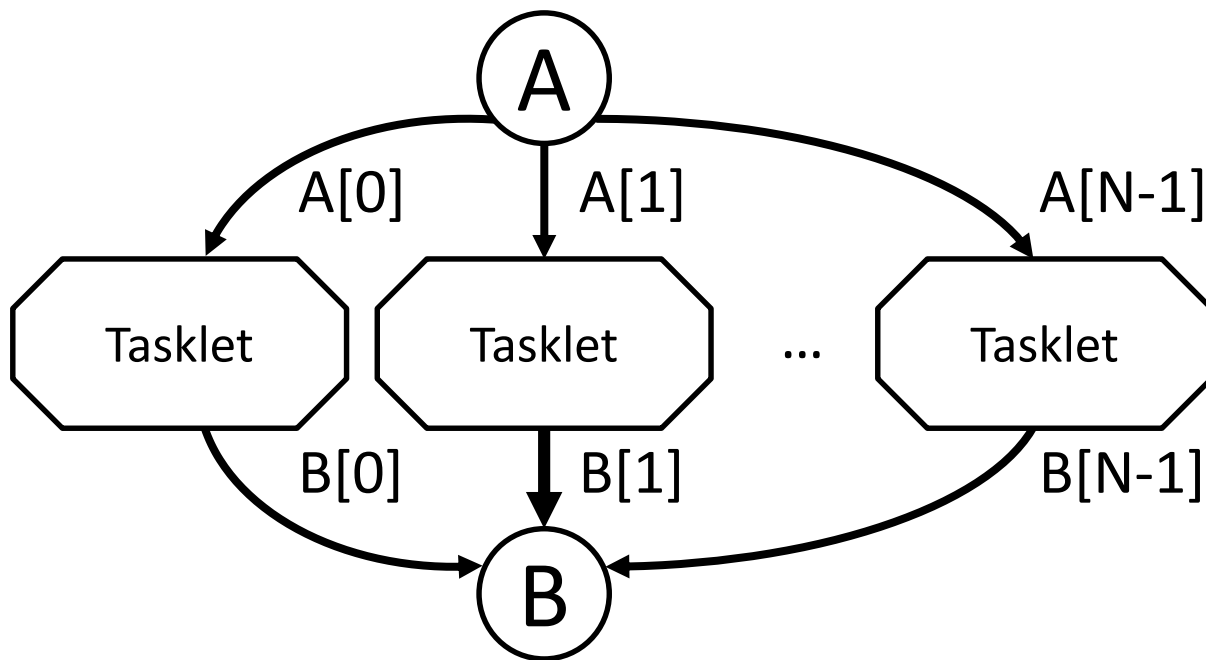
Dataflow Programming in DaCe



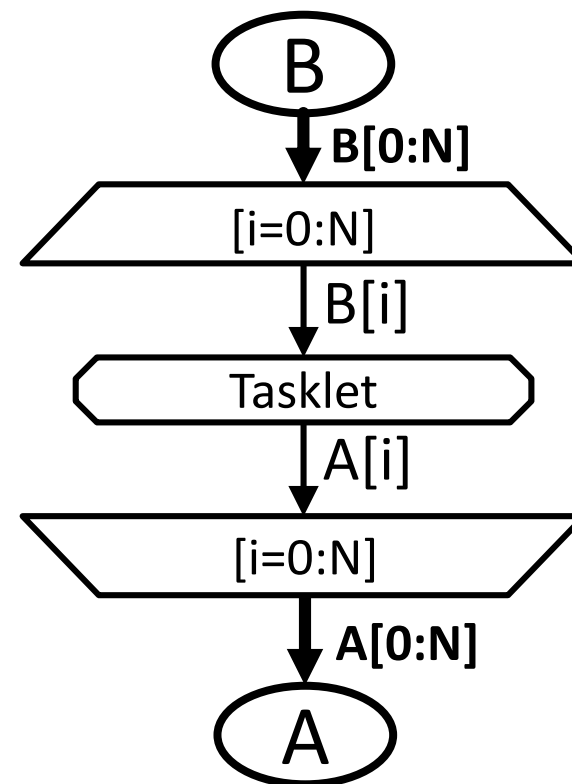
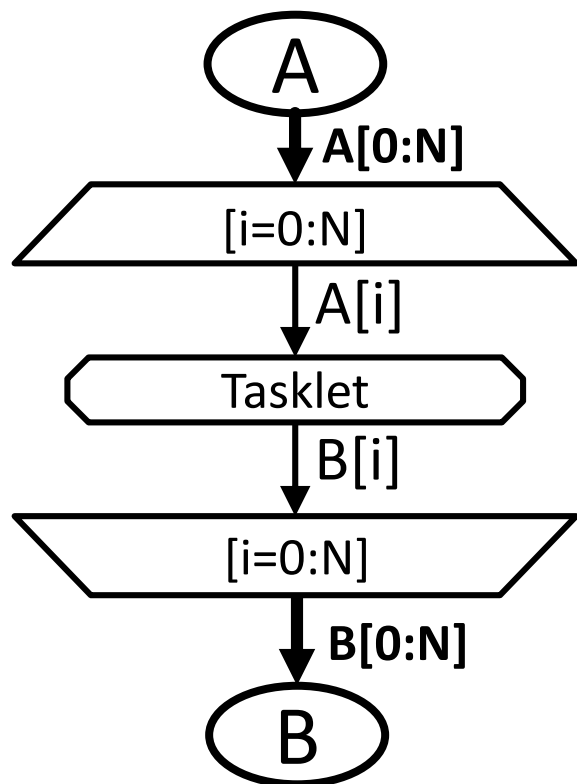
Parallel Dataflow Programming



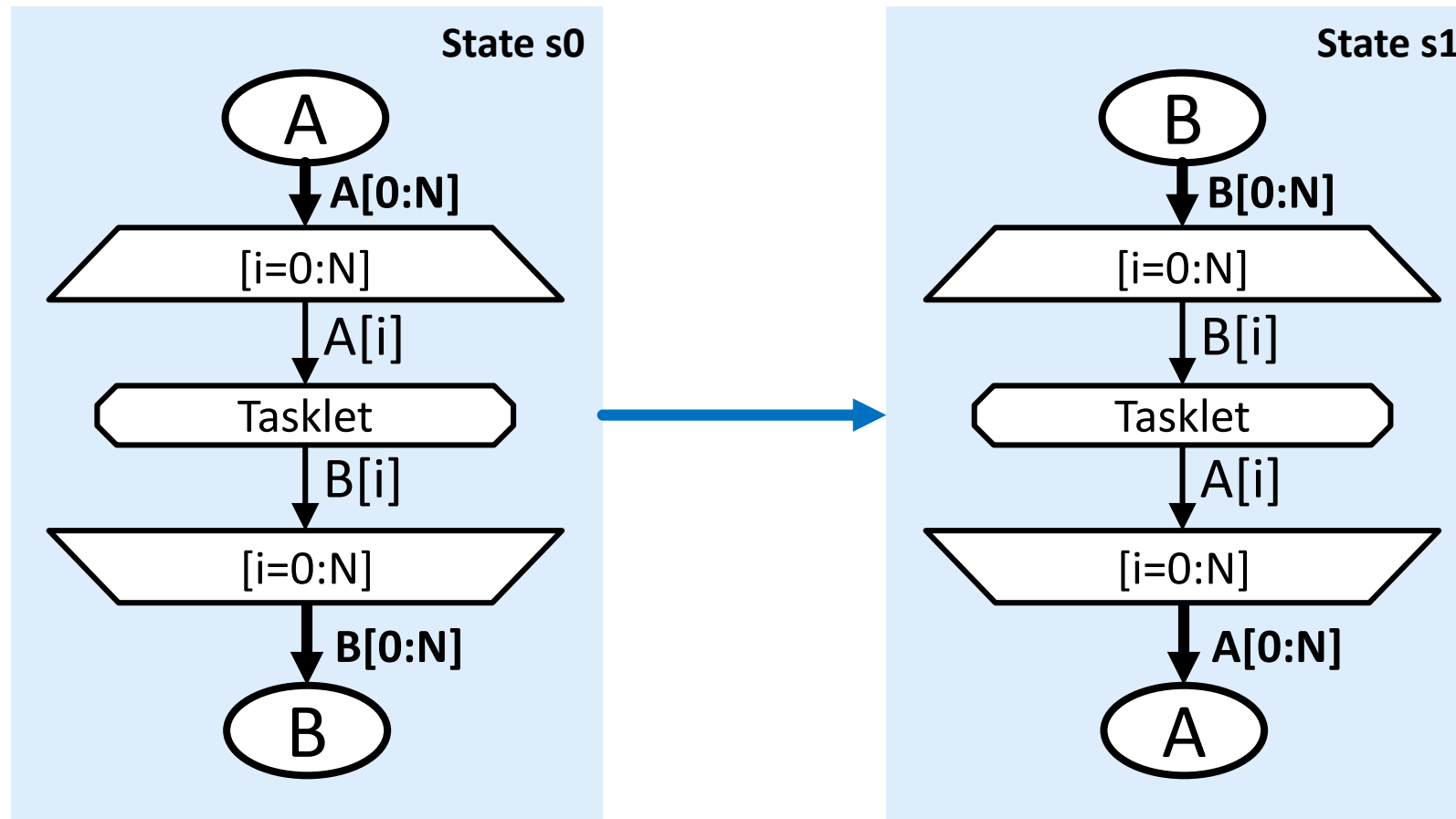
Parallel Dataflow Programming



Stateful Dataflow Parallel Programming in DaCe

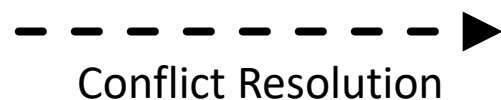
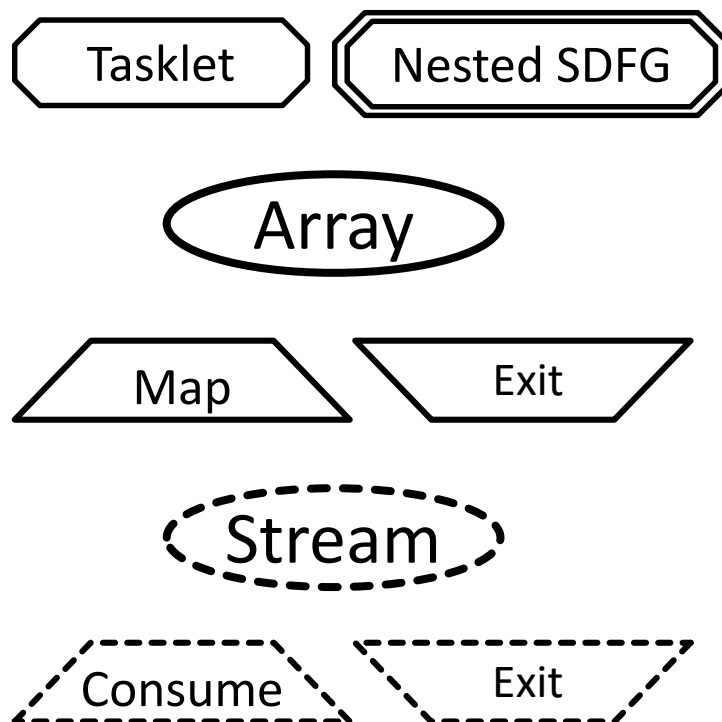


Stateful Dataflow Parallel Programming in DaCe



Meet the Nodes

State



State machine element

Fine-grained computation

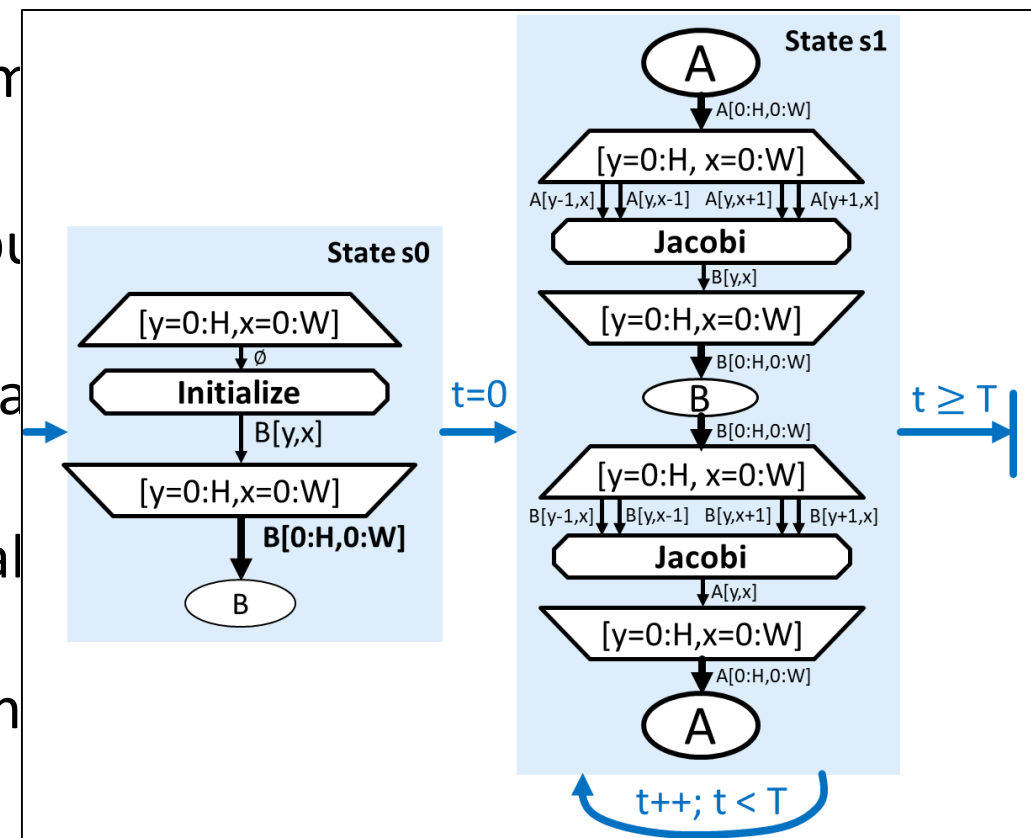
N-dimensional data

Parametric graph algorithm

Streaming data consumption

Dynamic mapping of computations on streams

Defines behavior during conflicting writes

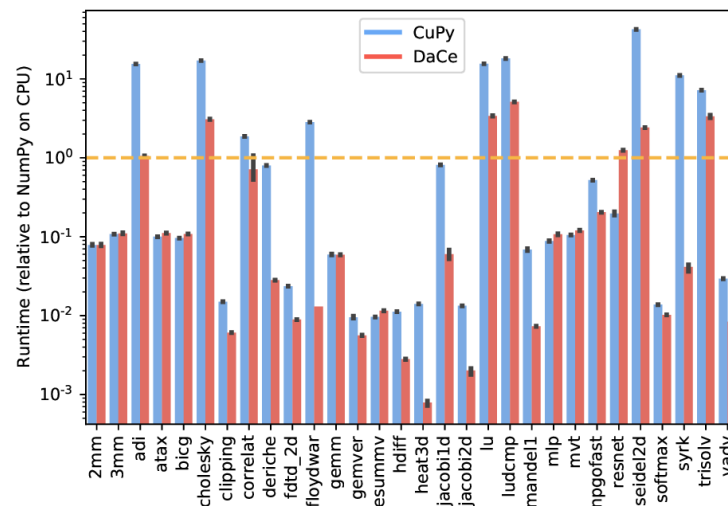


DaCe – Multi Backend Code Generation

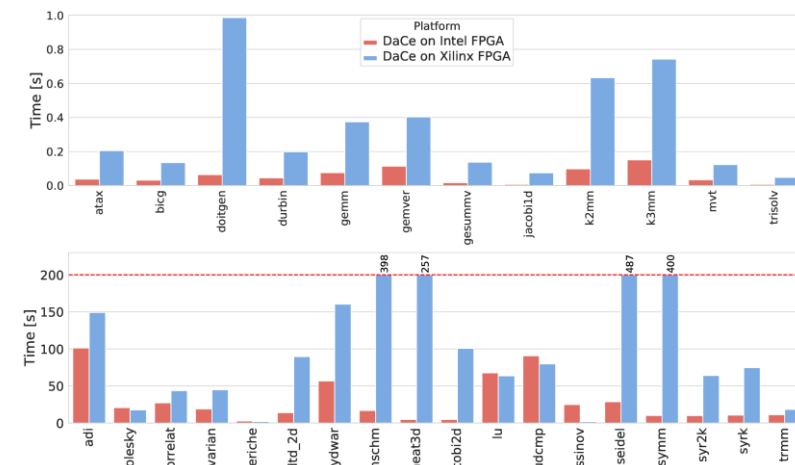
CPU Results

	DaCe	GCC	ICC	Numba	Pythran	NumPy
gramschm	↑9.1 ⁽²⁾	↑9.3	↑7.2 ⁽⁵⁴⁾	↑3.0	↑8.1	0.15 s
heat3d	↑454 ⁽⁴²⁾	↑24.0 ⁽⁹⁾	↑179 ⁽⁷⁾	↑50.1	↑2.3	50.35 s
jacobi1d	↑11.4 ⁽¹⁾	↑3.7 ⁽⁸⁾	↑3.9 ⁽¹³⁾	↓1.2	↑3.1	0.45 s
jacobi2d	↑56.2	↑7.1 ⁽²⁸⁾	↑58.6 ⁽⁷⁾	↑18.2	↑21.8	174.61 s
lu	↑3.6	↑2.9	↑4.5	↑2.2	↑4.8	13.02 s
ludcmp	↑3.7	↑5.4	↑5.4	↑2.3	↑4.9	13.7 s ⁽¹⁾
syr2k	↑10.4	↑10.5	↑287 ⁽¹⁴⁾	↑5.3	↑5.9	13.94 s
nussinov	↑612	↑1.1k	↑909 ⁽¹⁸⁾	↑420	↑871	20.37 s ⁽²⁾
seidel2d	↑185	↑158	↑78.1 ⁽²⁾	↑95.8	↑141	15.87 s ⁽¹⁾
symm	↑4.2	↑6.4	↑74.9 ⁽¹⁹⁾	↑11.5	↑4.2	9.88 s
syrk	↑10.0	↑9.3	↑233 ⁽²⁸⁾	↑3.0	↑5.5	6.58 s
gesummv	↑8.6 ⁽²⁾	↑1.8 ⁽⁶⁾	↑1.8 ⁽³⁾	↑7.5 ⁽⁵⁾	↑1.0	0.77 s
mvt	↑1.0	↓52.6 ⁽²³⁾	↓13.9 ⁽⁴⁾	↑1.0 ⁽¹⁾	↓1.0 ⁽¹⁾	45.22 ms ⁽¹⁾
gemver	↑19.1 ⁽¹³⁾	↑1.5 ⁽⁵⁾	↑2.5 ⁽¹⁾	↑2.1 ⁽³⁾	↑1.3 ⁽⁶⁾	0.79 s ⁽⁸⁾
atax	↓1.2 ⁽⁷⁾	↓12.0 ⁽²⁾	↓6.2 ⁽⁶⁾	↓1.0 ⁽⁵⁾	↓1.1 ⁽⁶⁾	73.37 ms ⁽⁵⁾
floydwar	↑116 ⁽²⁾	↑3.9 ⁽⁵⁾	↑2.0	↑57.3	↑7.3	84.0 s
fdtd_2d	↑170 ⁽⁹⁸⁾	↑3.7 ⁽¹¹⁾	↑41.3 ⁽²⁴⁾	↑4.1	↑1.3	7.4 s
durbin	↑2.2	↑5.0 ⁽¹⁾	↑5.9 ⁽¹⁰⁾	↓1.1 ⁽²⁾	↑3.2	0.65 s
doitgen	↑39.7 ⁽⁴⁶⁾	↓11.8	↓2.3 ⁽³⁾	↑1.1	1.0 ⁽³⁾	0.47 s ⁽⁸⁾
deriche	↑17.3	↑2.7 ⁽⁶⁾	↑55.1 ⁽⁹⁾	↑1.5	↑1.6	2.83 s
covarian	↑1.6 ⁽¹⁴⁾	↓21.4	↓5.0 ⁽⁸⁾	↑1.3 ⁽²⁾	↓20.4 ⁽³⁾	80.11 ms
correlat	↑1.6	↓27.8 ⁽²⁾	↓6.6 ⁽⁴⁾	↓1.1 ⁽⁹⁾	1.0 ⁽¹⁹⁾	61.55 ms ⁽¹⁹⁾
cholesky	↑12.1	↑3.7	↑3.4	↑10.3	↑15.2	7.05 s
bicg	↓1.1 ⁽⁸⁾	↓17.2 ⁽²⁾	↓17.5 ⁽¹⁾	↓1.1 ⁽⁷⁾	↓1.1 ⁽⁶⁾	75.92 ms ⁽⁹⁾
trisolv	↑2.3 ⁽¹⁴⁾	↓1.6 ⁽²⁾	↓1.8 ⁽⁴⁾	↑1.9 ⁽²³⁾	↑1.3	0.13 s ⁽¹⁴⁾
adi	↑16.5	↑9.1	↑13.7 ⁽⁶⁾	↑8.2	↑9.4	0.9 s
3mm	↑1.7 ⁽³⁾	↓1.7k ⁽¹¹⁾	↓1.1 ⁽²⁾	↑1.6 ⁽⁶⁾	↑1.2 ⁽¹³⁾	0.46 s ⁽²¹⁾
2mm	↑2.4 ⁽²⁾	↓942 ⁽²²⁾	↓17.1	↑1.5 ⁽¹⁰⁾	↑1.3 ⁽¹⁶⁾	0.42 s ⁽¹²⁾
gemm	↑2.3 ⁽³⁾	↓133 ⁽²²⁾	↓4.8 ⁽¹¹⁾	↑1.4 ⁽⁷⁾	↑1.4	78.72 ms
trmm	↑36.2	↑4.1	↑50.9 ⁽¹⁸⁾	↑15.3 ⁽⁴⁾	↑4.2	4.12 s ⁽¹⁾

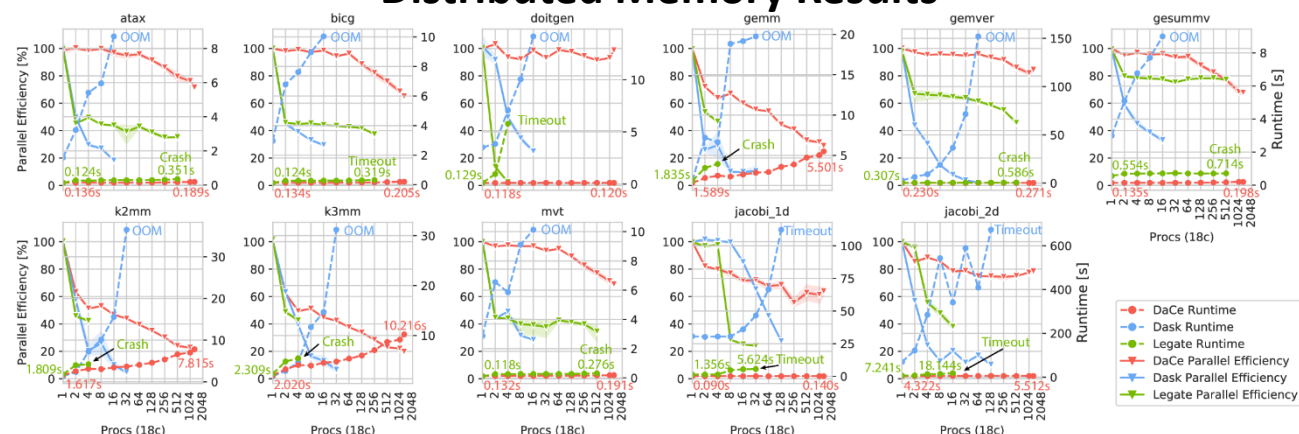
GPU Results



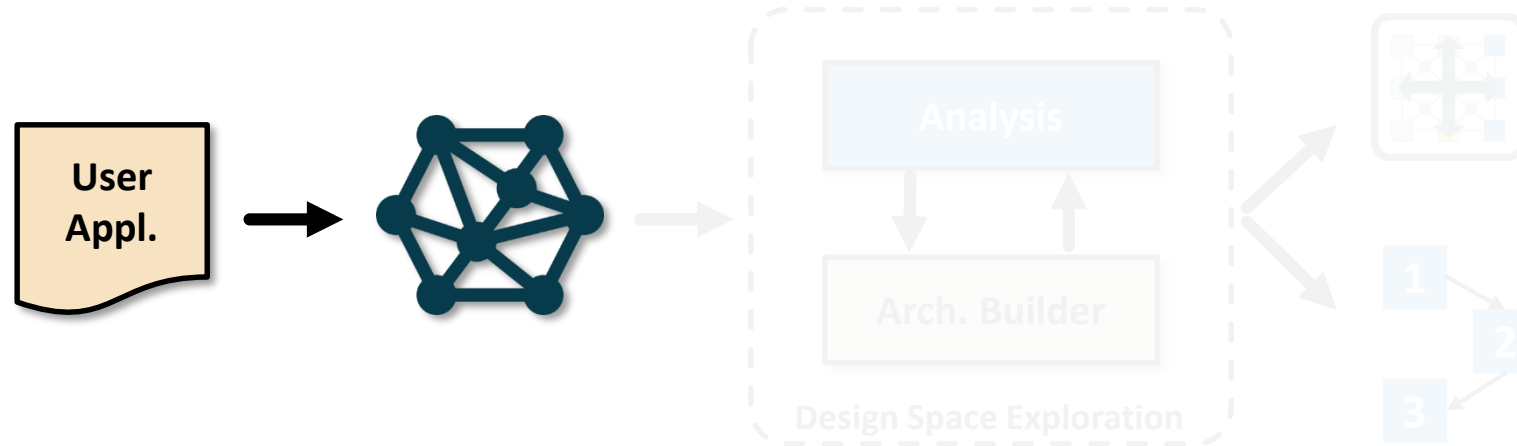
FPGA Results



Distributed Memory Results



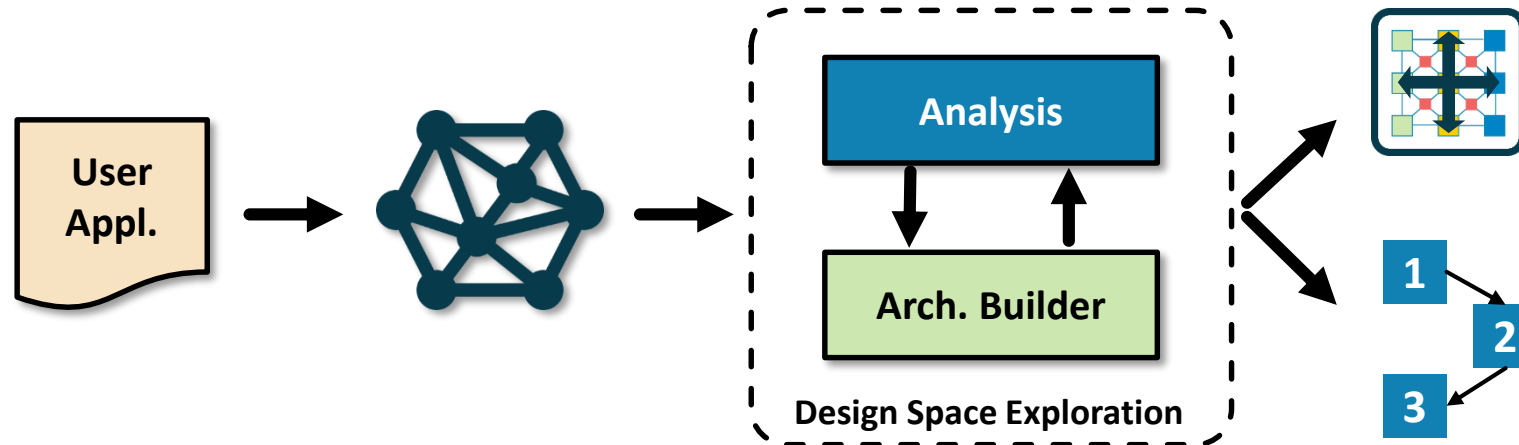
First steps



In this project we want to utilize DaCe to:

- exploit **its Data-Centric Intermediate Representation**
- (later on) use it as **“frontend”**, to let user express her own application using a high-level formalism

First steps



Understand how we can efficiently **map an application on a given architecture**:

- we want to consider **heterogeneity**
- we want to deal with architecture specific features, such as the presence of a **fast on-chip interconnect**

Your inputs regarding typical workloads/applications/use cases will be important



Thank you!

spcl.inf.ethz.ch

@spcl_eth

ETH zürich