

# Design Principles of Dynamic Resource Management in High-Performance Parallel Programming Models

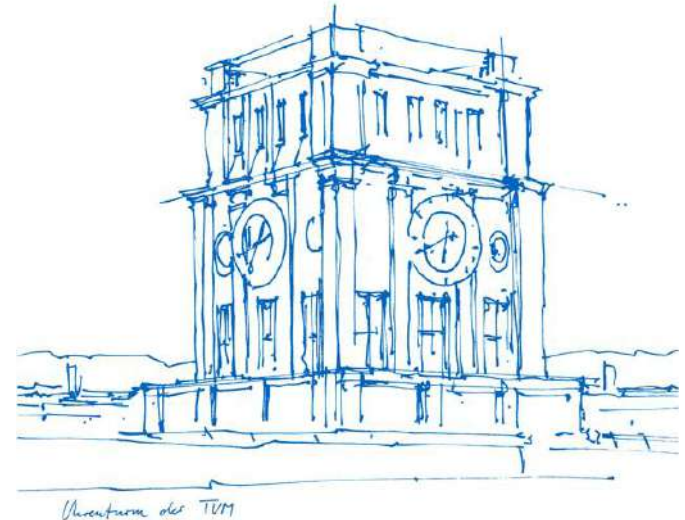
Dominik Huber<sup>1</sup>, Martin Schreiber<sup>2</sup>, Martin Schulz<sup>1</sup>, Howard Pritchard<sup>3</sup>, Daniel Holmes<sup>4</sup>

<sup>1</sup>Technical University of Munich, <sup>2</sup>Université Grenoble Alpes

<sup>3</sup>Los Alamos National Laboratory, <sup>4</sup>Intel Corporation

**4th EuroHPC Workshop on Dynamic Resources in HPC**

Dresden, 25.08.2025



# OUTLINE

## **The Need for Common Interfaces**

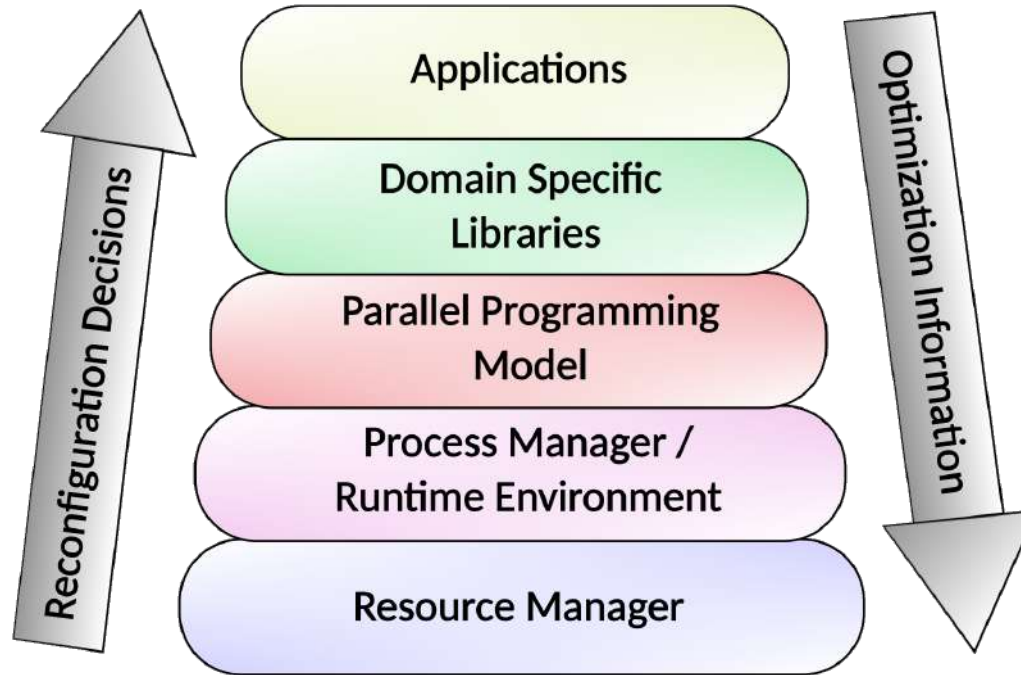
Design Principles

Minimal Interfaces

Prototype & Early Experiences



# Dynamic Resources in the HPC Software Stack



# The Need for Common Interfaces

## State-of-the-art

- Specialized solutions
- Not portable/generalizable
- No standardization



## Future?

- Common interfaces
- Specialization on top
- Standardization



# OUTLINE

The Need for Common Interfaces

**Design Principles**

Minimal Interfaces

Prototype & Early Experiences

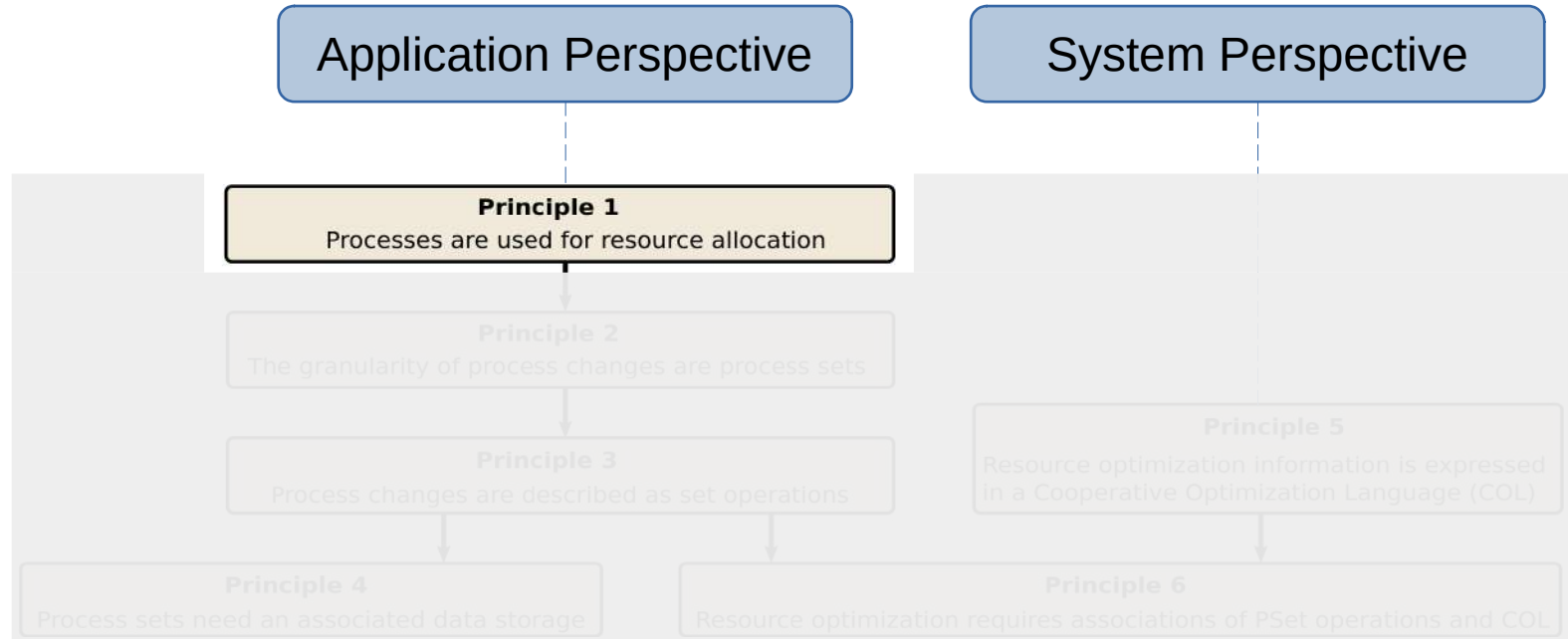


# Principle-Driven Design

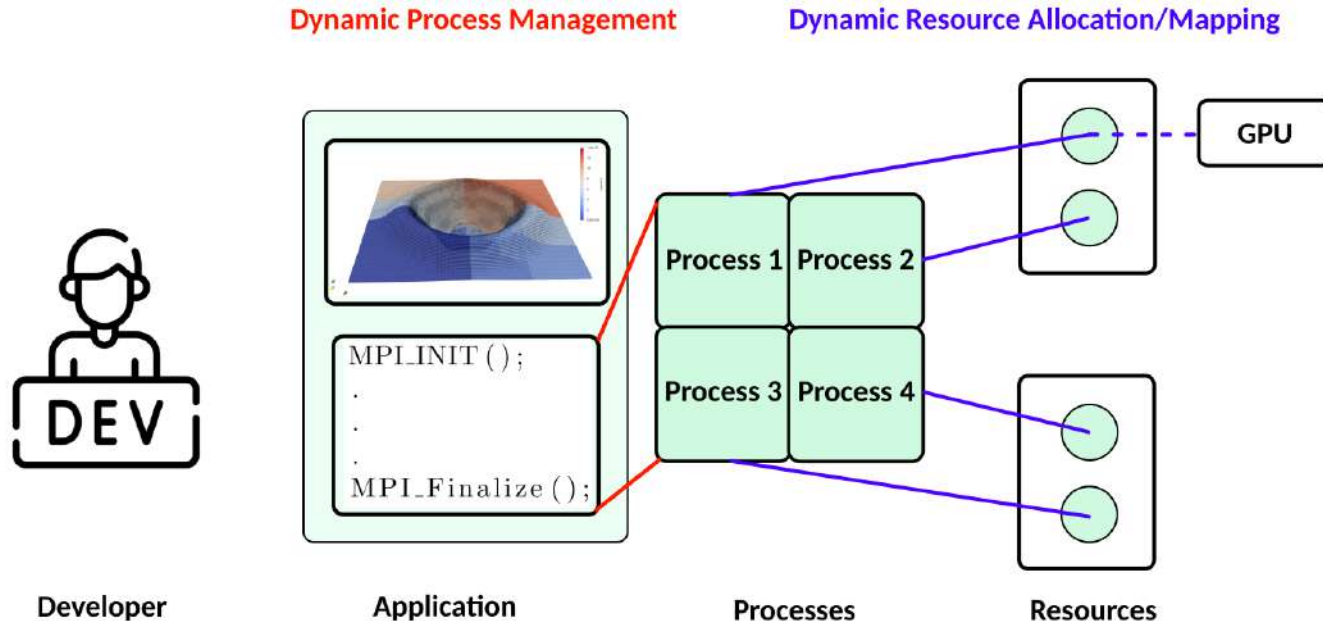
Common design principles for DRM mechanisms in the middleware stack



# Principle-Driven Design

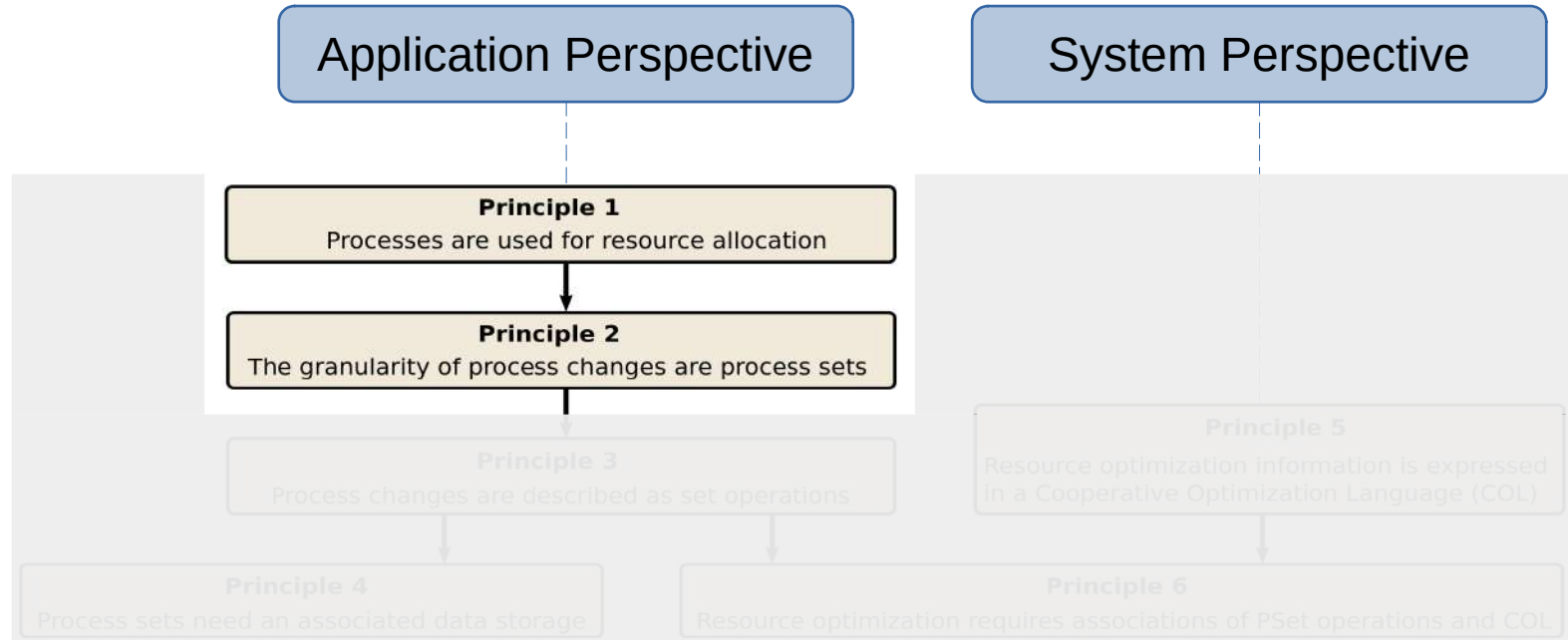


# Principle 1: Processes are used for resource allocation





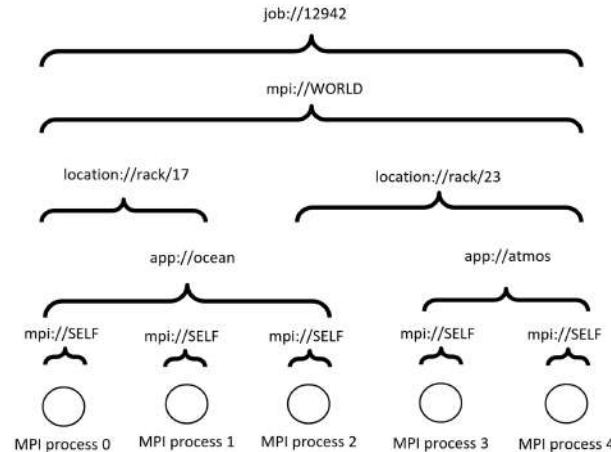
# Principle-Driven Design



## Principle 2: The granularity of process changes are PSets

**Definition “Process Set” (PSet):** A PSet is a set of processes where both, the set and its members are uniquely identifiable. A “**0-Pset**” is an empty Pset which always exists.

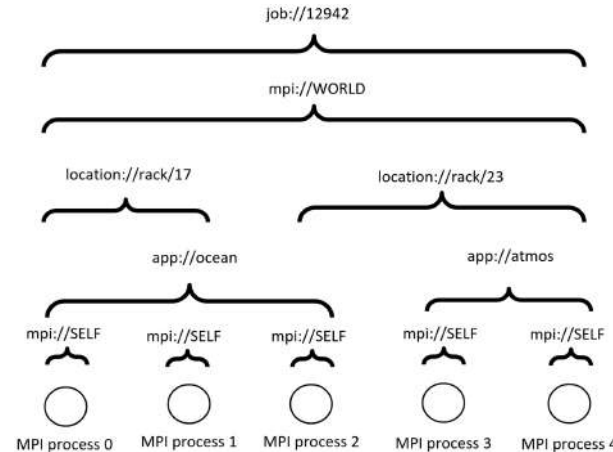
**Definition “Process Changes”:** A process change is a change of association between processes/PSets and HPC resources.



**Example of PSets in MPI** (MPI: A Message-Passing Interface Standard Version 4.0. p. 504)

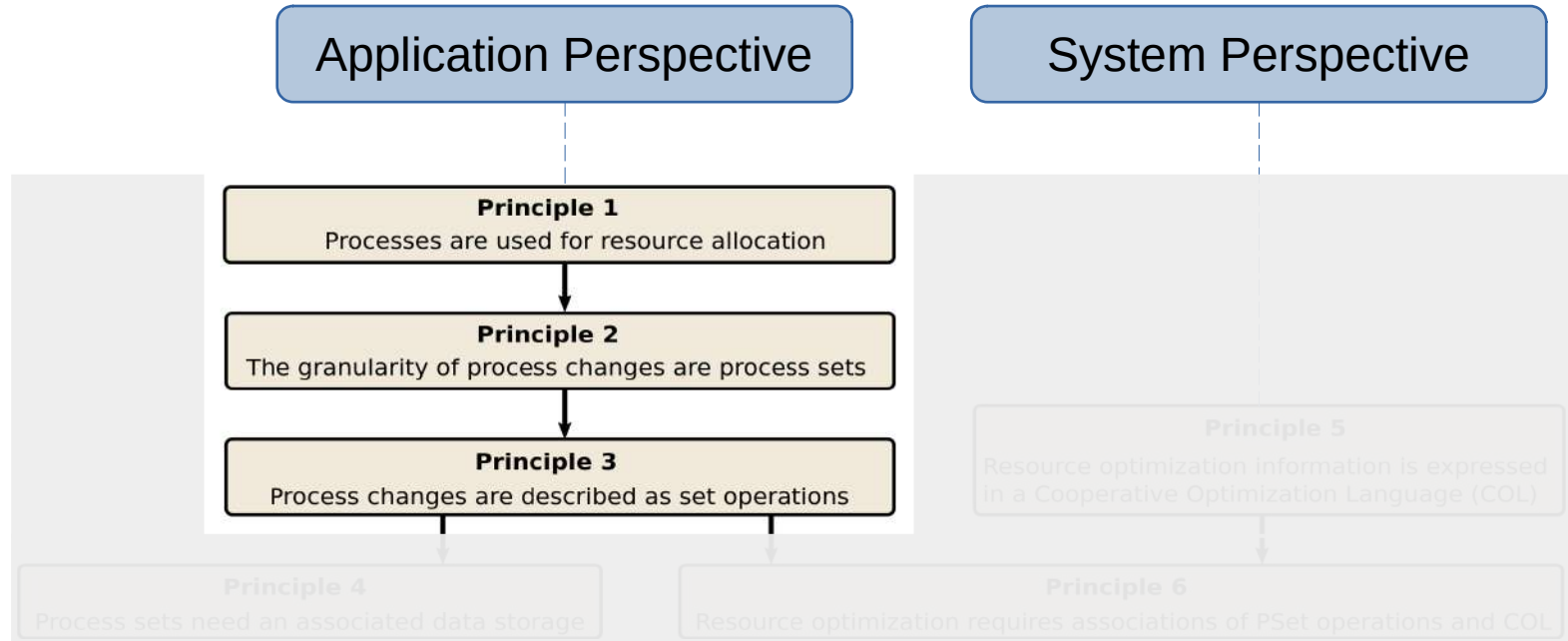
# Principle 2: The granularity of process changes are PSets

1. PSets provide **flexible granularity**
2. Unique PSet labels allow **reference without listing members**
3. Unique Proc labels are required for **orderability and set operations**



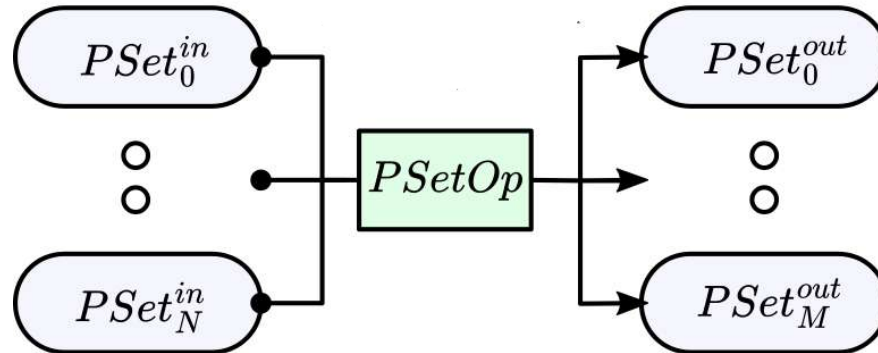
Example of PSets in MPI (MPI: A Message-Passing Interface Standard Version 4.0. p. 504)

# Principle-Driven Design



## Principle 3: Process changes are described as set operations

- **Definition “PSet Operation” (PSetOp):** A PSetOp is an operation on a list of **input PSets** that produces a list of **output PSets** based on a **well-defined rule**.



## Principle 3: Process changes are described as set operations

- **Definition “PSet Operation” (PSetOp):** A PSetOp is an operation on a list of **input PSets** that produces a list of **output PSets** based on a **well-defined rule**.

ADD:  $P_{in}^0 \rightarrow P_{out}^0, P_{in}^0 \cap P_{out}^0 = \emptyset$

Expresses launch of a task with a new set of processes.

SUB:  $P_{in}^0 \rightarrow \emptyset$

Expresses termination of a task executed by a set of processes.

SPLIT:  $P_{in}^0 \rightarrow P_{out}^k, \bigcap_i P_{out}^i = \emptyset \wedge \bigcup_i P_{out}^i = P_{in}^0$

Expresses splitting of a task into subtasks.

UNION:  $P_{in}^i \rightarrow P_{out}^0, \bigcup_i P_{in}^i = P_{out}^0$ ,

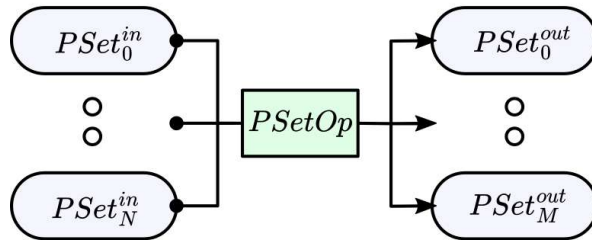
Expresses merging subtasks.

GROINK:  $P_{in}^0 \rightarrow \{P_{out}^0, P_{out}^1, P_{out}^2\}, P_{out}^1 \subset P_{out}^0 \wedge P_{out}^2 = P_{in}^0 \setminus P_{out}^0 \cup P_{out}^1$ .

Expresses adding/removing processes to/from a task executed by a PSet.

# Principle 3: Process changes are described as set operations

- **Rationale:**
  - Relates process changes to existing PSets
  - Generic abstraction for different scenarios



ADD:  $P_{in}^0 \rightarrow P_{out}^0, P_{in}^0 \cap P_{out}^0 = \emptyset$   
 Expresses launch of a task with a new set of processes.

SUB:  $P_{in}^0 \rightarrow \emptyset$   
 Expresses termination of a task executed by a set of processes.

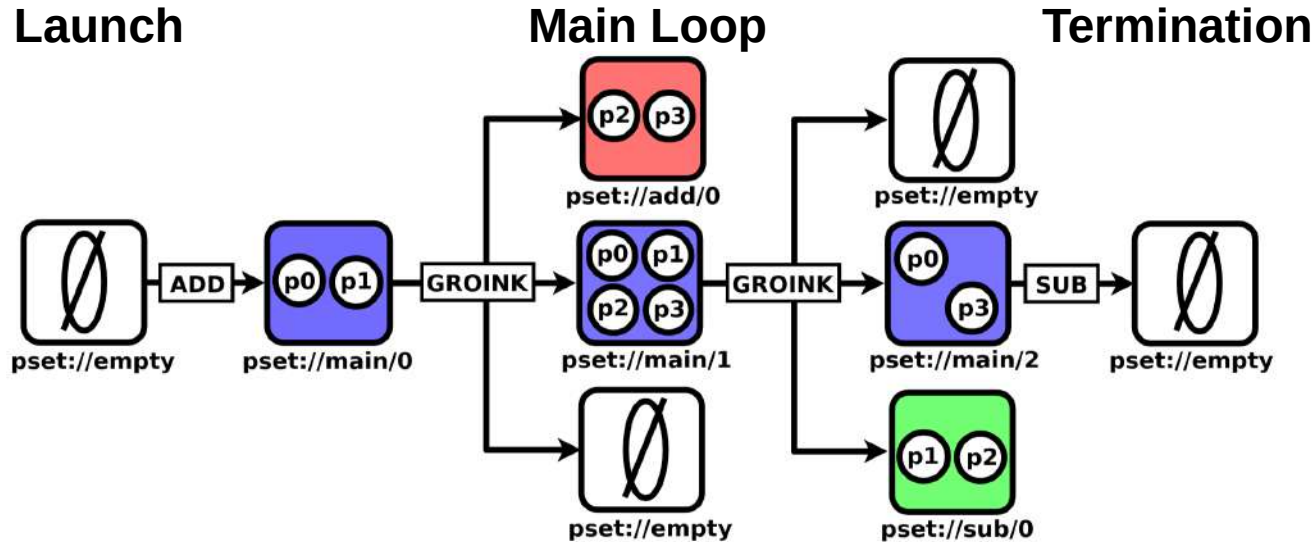
SPLIT:  $P_{in}^0 \rightarrow P_{out}^k, \bigcap_i P_{out}^i = \emptyset \wedge \bigcup_i P_{out}^i = P_{in}^0$   
 Expresses splitting of a task into subtasks.

UNION:  $P_{in}^i \rightarrow P_{out}^0, \bigcup_i P_{in}^i = P_{out}^0$   
 Expresses merging subtasks.

GROINK:  $P_{in}^0 \rightarrow \{P_{out}^0, P_{out}^1, P_{out}^2\}, P_{out}^1 \subset P_{out}^0 \wedge P_{out}^2 = P_{in}^0 \setminus P_{out}^0 \cup P_{out}^1$   
 Expresses adding/removing processes to/from a task executed by a PSet.

# Principle 3: Process changes are described as set operations

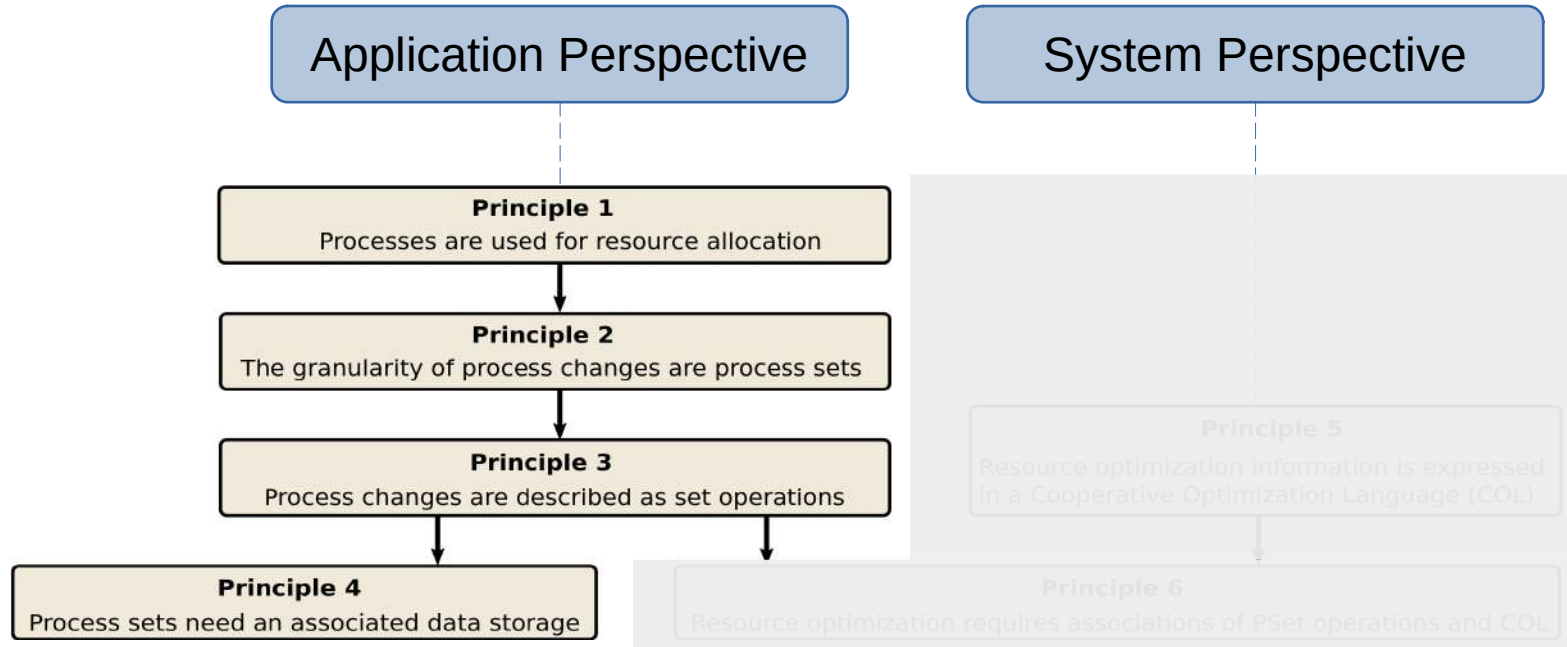
- **Example:** Resizing MPI\_COMM\_WORLD



D. Huber, S. Iserte, M. Schreiber, A. J. Peña and M. Schulz, "Bridging the Gap Between Genericity and Programmability of Dynamic Resources in HPC," ISC High Performance 2025 Research Paper Proceedings (40th International Conference), Hamburg, Germany, 2025, pp. 1-11.



# Principle-Driven Design

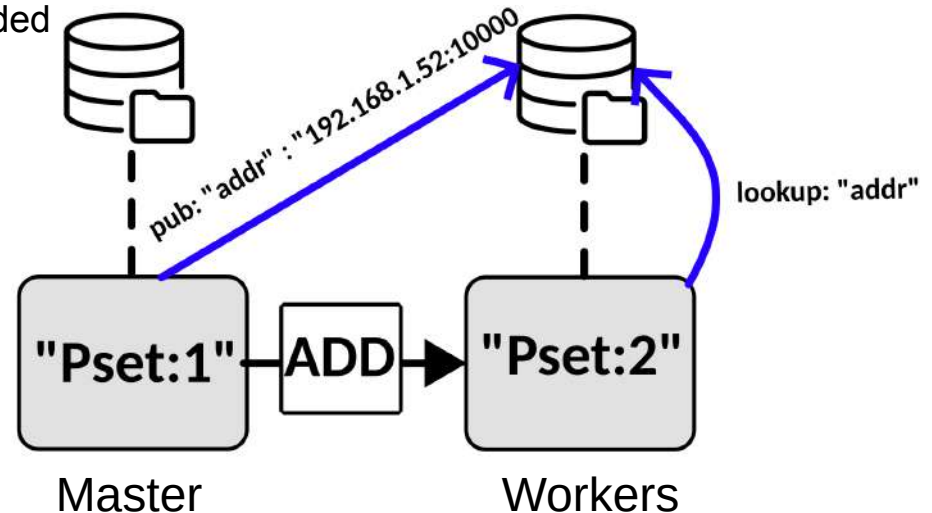


## Principle 4: Process Sets need an associated data storage

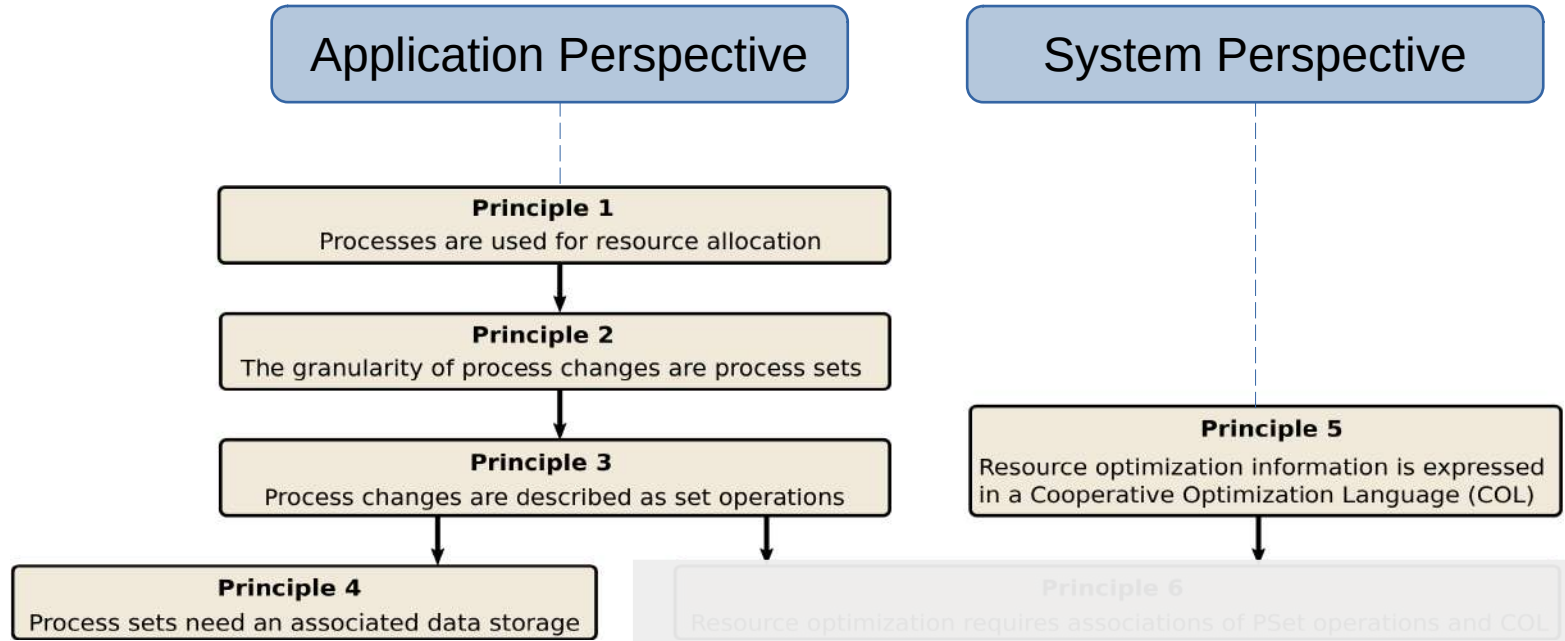
### Dynamic Process Environment:

- Possibly no direct communication channels between PSets/processes
- Necessary information not always known a-priori
- Generic way of exchanging „control“ data needed

**Example:** Master-Worker



# Principle-Driven Design



## Principle 5: Optimization Information is expressed in a COL

**Definition “Cooperative Optimization Language (COL)”:** A Cooperative Optimization Language expresses optimization information for DRM in a **cooperative**, **local**, and **language-based** way.

## Principle 5: Optimization Information is expressed in a COL

**Definition “Cooperative Optimization Language (COL)”**: A Cooperative Optimization Language expresses optimization information for DRM in a **cooperative**, **local**, and **language-based** way.

### Cooperative

- Actors contribute local information toward a global goal
- Global optimization supports local resource needs

## Principle 5: Optimization Information is expressed in a COL

**Definition “Cooperative Optimization Language (COL)”**: A Cooperative Optimization Language expresses optimization information for DRM in a **cooperative**, **local**, and **language-based** way.

### Cooperative

- Actors contribute local information toward a global goal
- Global optimization supports local resource needs

### Local

- Provide optimization information solely for themselves
- No reliance on global system knowledge
- Resource request are „implicit“ through provided optimization information

## Principle 5: Optimization Information is expressed in a COL

**Definition “Cooperative Optimization Language (COL)”**: A Cooperative Optimization Language expresses optimization information for DRM in a **cooperative**, **local**, and **language-based** way.

### Cooperative

- Actors contribute local information toward a global goal
- Global optimization supports local resource needs

### Language-based

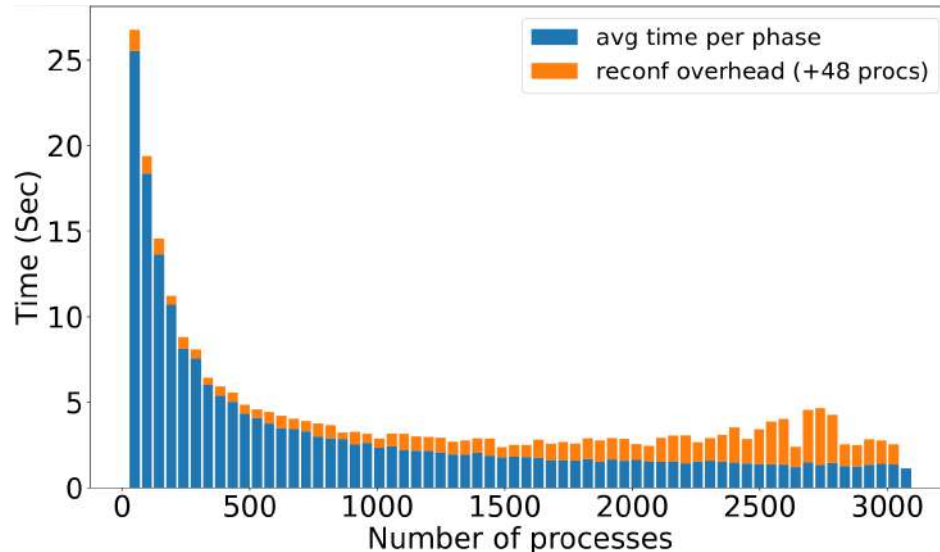
- A domain-specific language
- Intermediate Representation

### Local

- Provide optimization information solely for themselves
- No reliance on global system knowledge
- Resource request are „implicit“ through provided optimization information

# Principle 5: Optimization Information is expressed in a COL

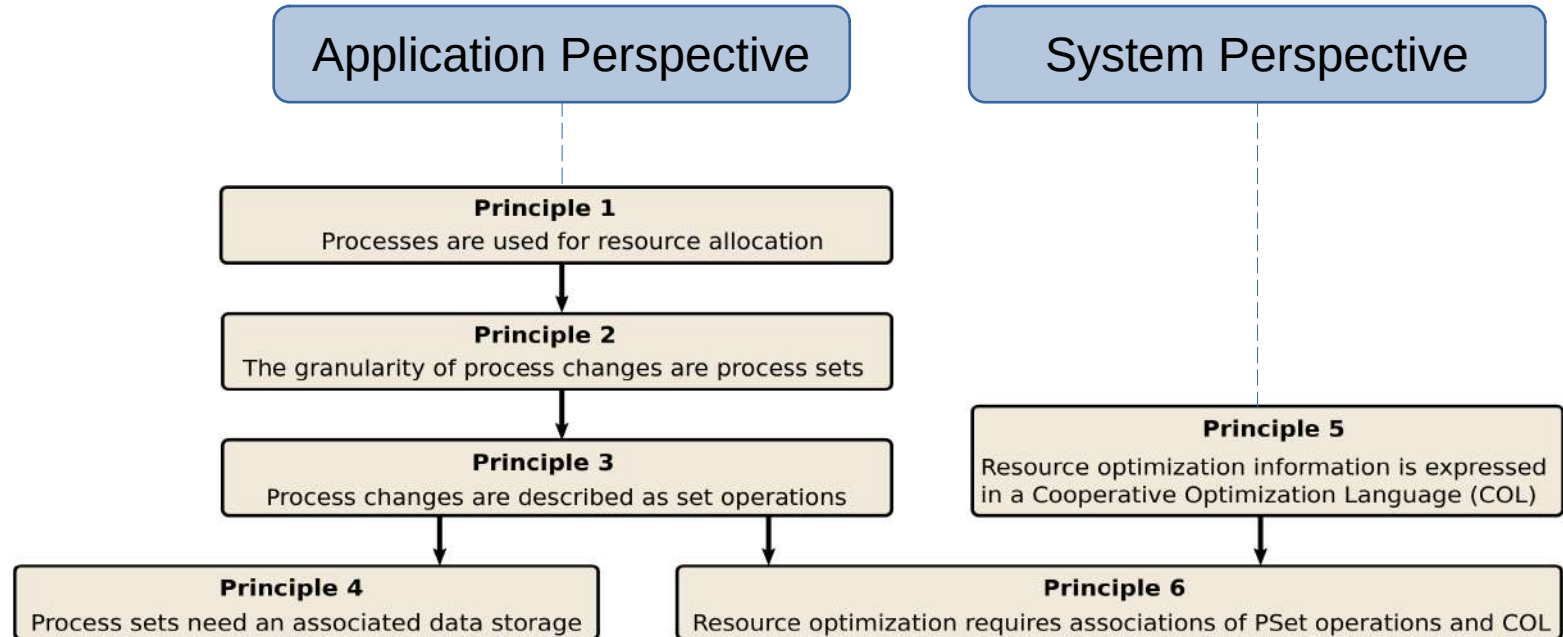
## Example: Simple optimization information for iterative workload



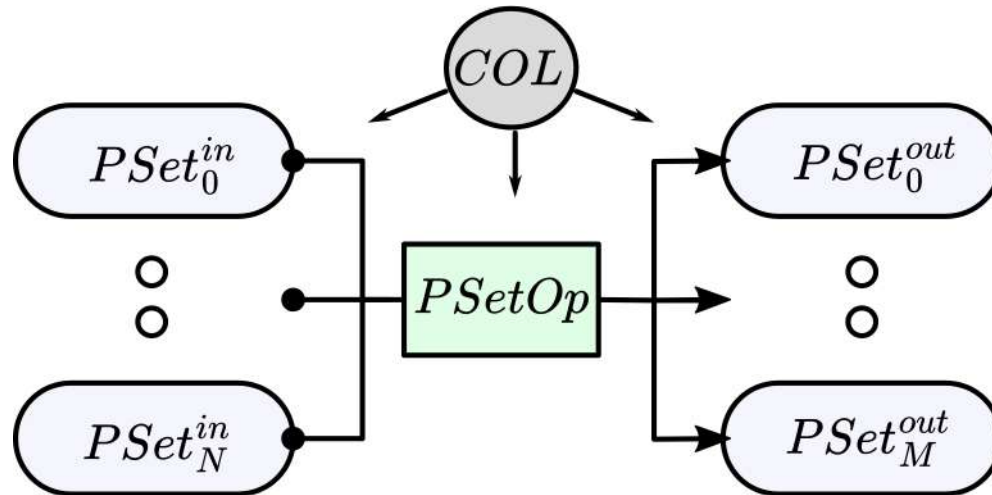
```
{
  opt_info : {
    perf_model : amdahl
    perf_params : t_s=1,t_p=300
    oh_model : 1.5*N + 4
    min_ram_per_proc : 10
    procs_per_node : 48
    min_procs : 96
    max_procs : 9600
  }
}
```



# Principle-Driven Design



## Principle 6: Association between Pset Operations and COL



# OUTLINE

The Need for Common Interfaces

Design Principles

**Minimal Interfaces**

Prototype & Early Experiences



## Minimal Interfaces: PSet Operations

```
// Indicates support for a PSetOp
async (op, in_sets, col_obj) -> out_sets
    op:          The type of the PSetOp
    in_sets:      List of input PSet names
    col_obj:      Info expressed by the COL
    out_sets:     List of output PSet names
```

**Listing 3.** Minimal interface for indicating support for a PSetOp.

## Minimal Interfaces: PSet Dictionary

```
// Publishes key-value pairs
(pset_name, key_vals) -> None
    pset_name:  ID for data storage
    key_vals:   Key-values to publish

// Retrieves key-value pairs
(pset_name, keys) -> key_vals
    pset_name:  ID for data storage
    keys:       Keys to be looked up
    key_vals    Retrieved key-values
```

**Listing 4.** Minimal interface for accessing the data storage associated with a PSet

# OUTLINE

The Need for Common Interfaces

Design Principles

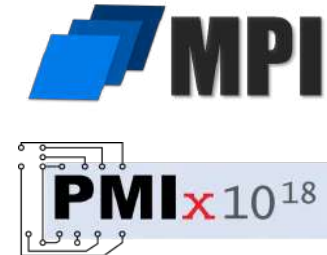
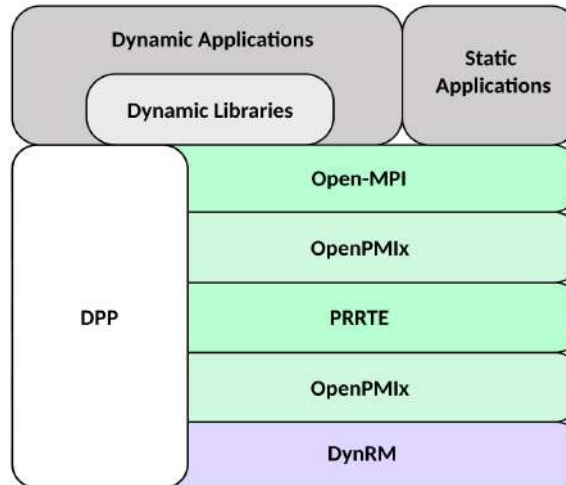
Minimal Interfaces

**Prototype & Early Experiences**



# Prototype and Early Experiences

<b>PetSc [1]</b> Parallel numerical software library	<b>DynLAIK [13]</b> Data distribution for dynamicity and fault tolerance	<b>LibPFASST [3]</b> Parallel-in-time integration library
<b>P4est [2]</b> Adaptive Mesh Refinement library	<b>XBraid [4]</b> Multigrid parallel-in-time integration	<b>MPDATA3D [12]</b> 3D semi-Lagrangian multiscale fluid solver
<b>DMR-API [6]</b> Dynamic resources API for computational kernels	<b>AMT-GLB [11]</b> Asynchronous Many-Task Runtime System	<b>DPPIInSitu [10]</b> Library for dynamic, asynchronous in-situ techniques



# Summary and Future Work

- **Goal: Common, standardized Dynamic Resource Management Interfaces**
- **Approach:**
  - **Generic Design Principles** for parallel programming models and middleware
  - **Prototype** and collaborative development environment
- **Future Work:**
  - **Standardization:** COL, MPI, PMIx, ...
  - **Integration** with production resource managers

## Contact

[domi.huber@tum.de](mailto:domi.huber@tum.de)

