

Techniques for Dynamic Workload Partitioning in High-Performance Heterogeneous Computing Platforms

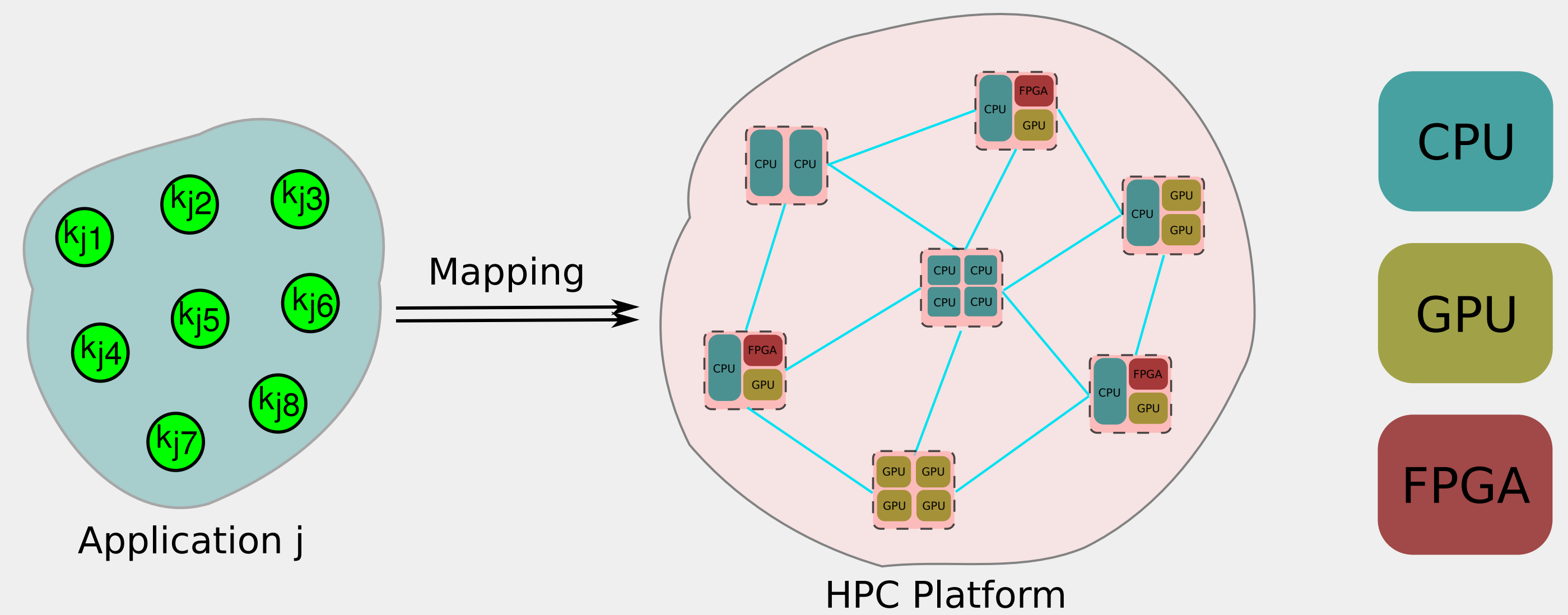


THALES

Wilder Lopes
wilderlopes@gmail.com

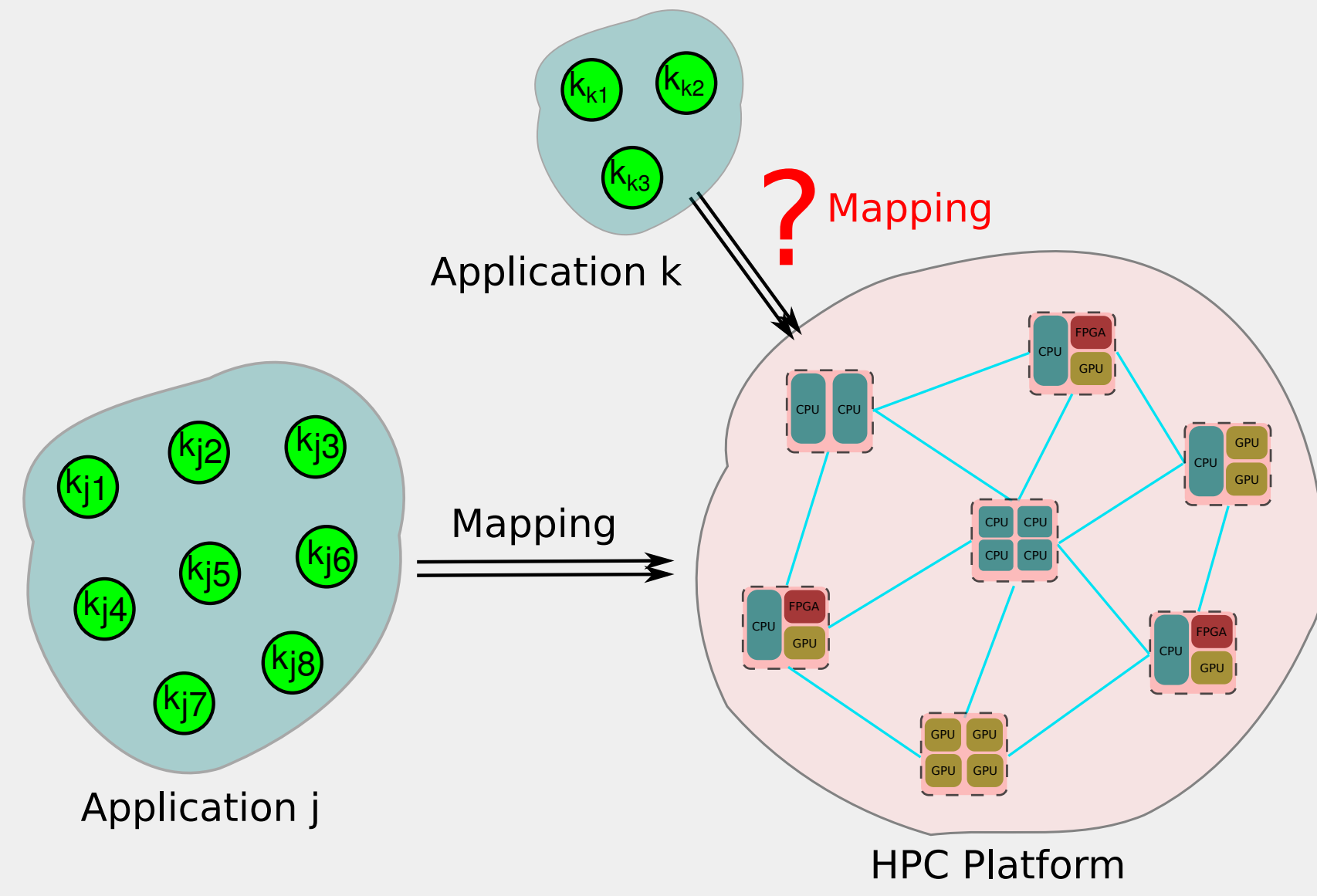
Introduction

- The proper allocation/mapping of data-parallel processes (kernels) in a High-Performance Computing (HPC) platform is crucial to exploit the system full potential
- Dynamic Heterogeneous Platforms: type and amount of available devices **may change at runtime**
- Applications requirements (Quality of Service – QoS): **also may change at runtime**

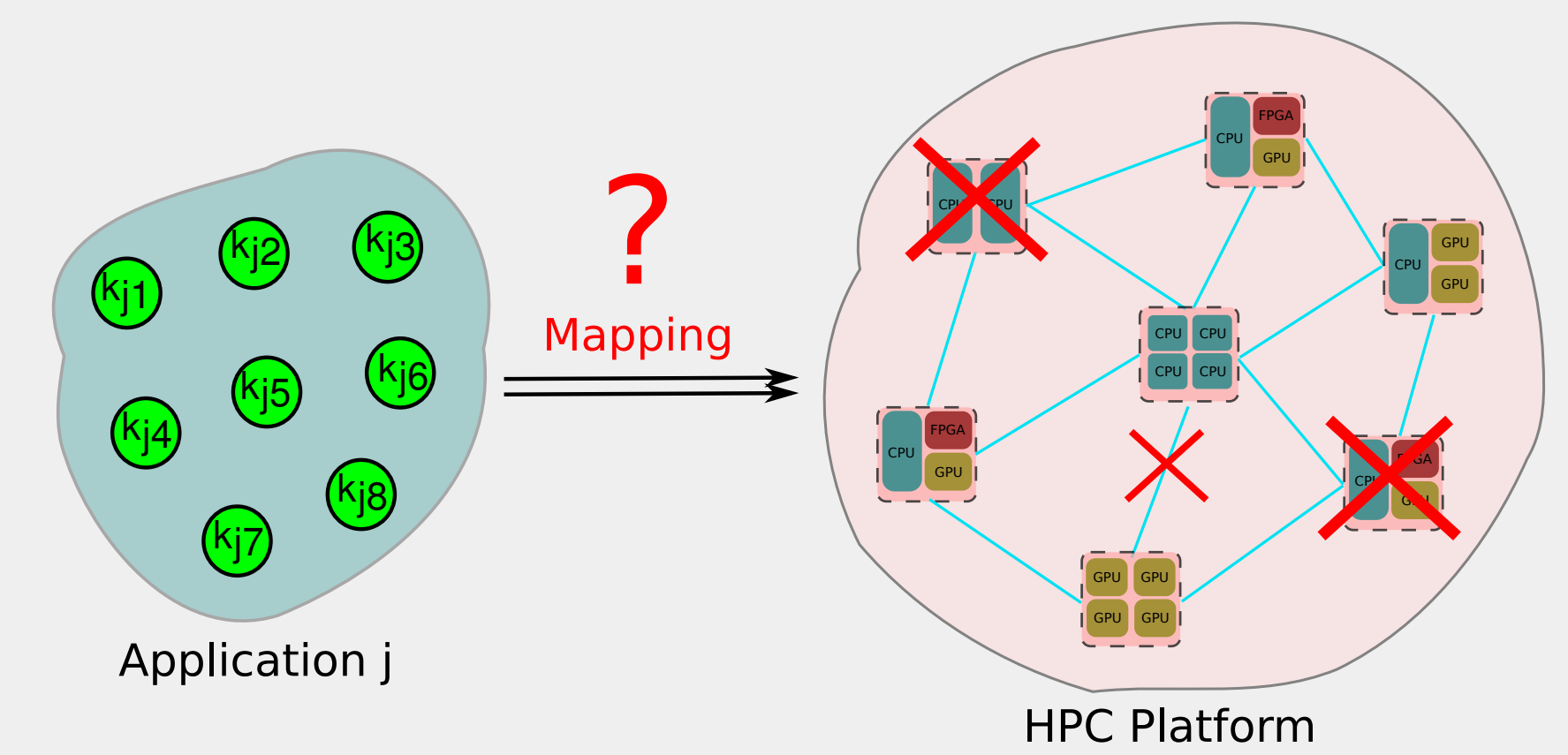


GOAL

Develop a system manager able to sense and react at runtime to variations in the High-Performance Heterogeneous Computing platform as well as in the QoS requirements



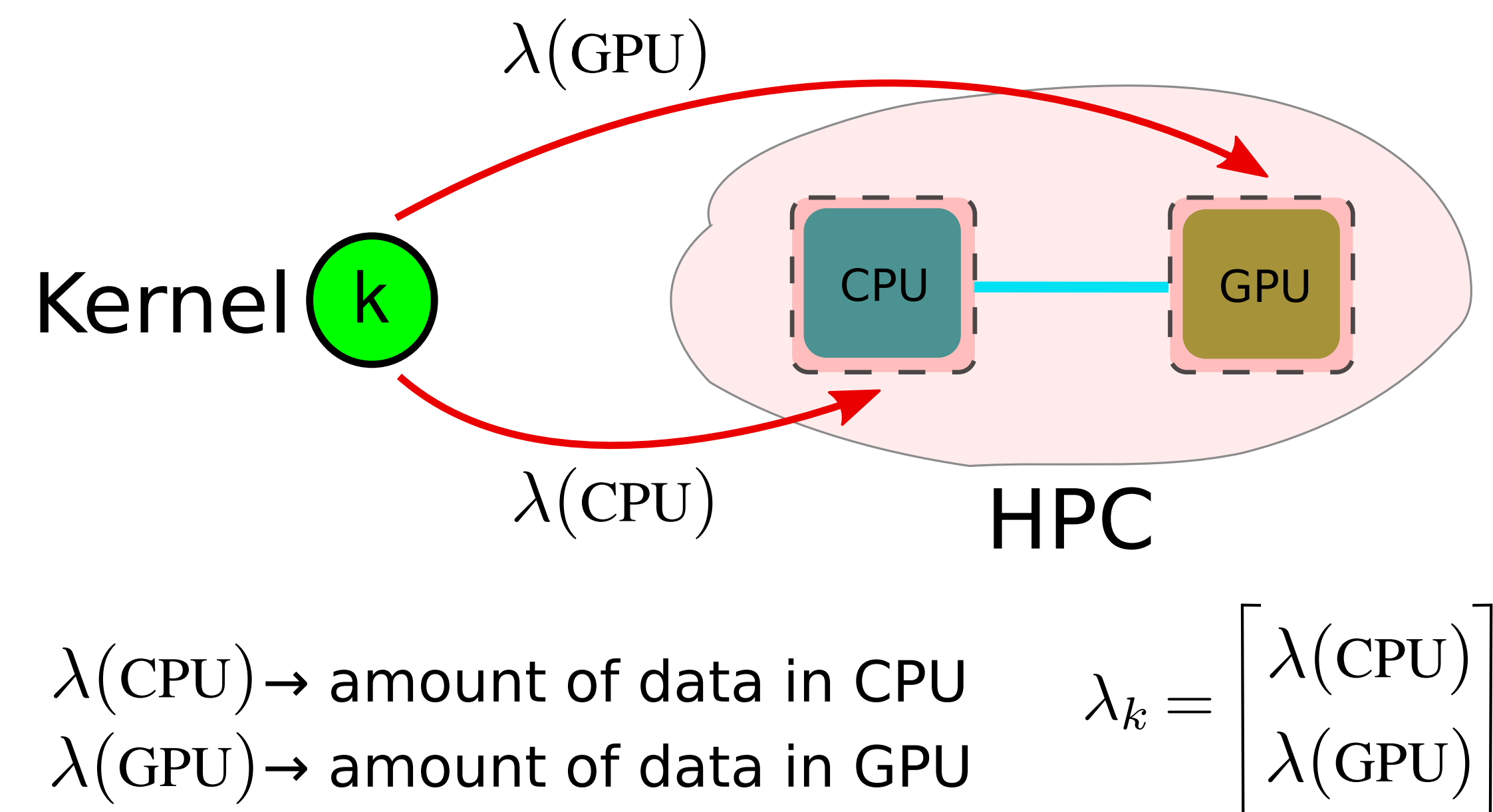
Scenario 1: new application needs to be mapped



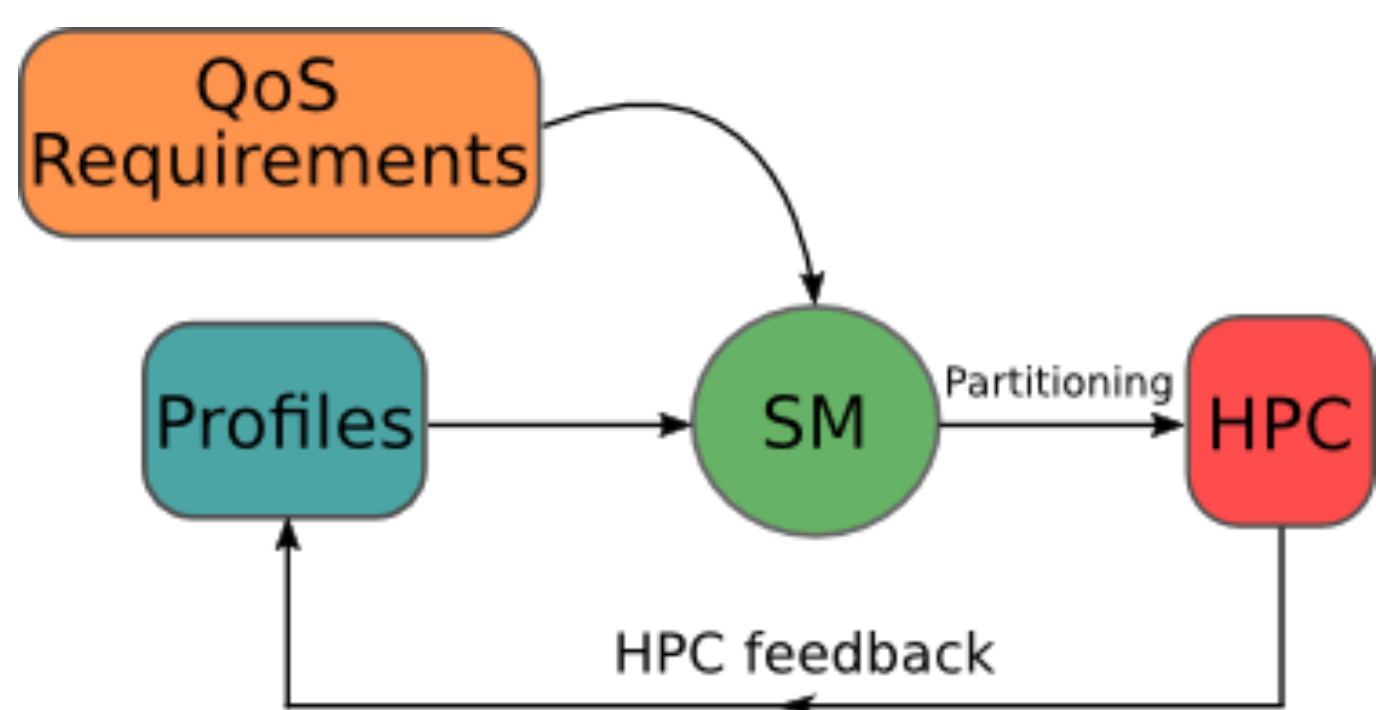
Scenario 2: resource/network failures

Problem Formulation

- Simple scenario:** single-kernel data-parallel application → preliminary study
- Only two devices: **CPU and GPU**
- Workload partitioning ↔ Data partitioning
- Estimate the **amount of data** to be allocated to each device
- Definition: Best partitioning** is the one that enforces minimization of execution time



Design of the System Manager (SM)



- It is assumed a **performance profile** of the kernel in each device is available → Matrix R
- HPC feedback:** measurement of devices performances in order to update profiles in R
- QoS requirements:** application dependent
- SM uses those to estimate the workload partitioning that minimizes kernel execution time

Mathematical Representation

error e_k : Measure of how well the resources can be combined to achieve the QoS requirements

$$e_k = \left(d - \sum_j r_j \lambda(j) \right) = (d - R\lambda_k)$$

$d \rightarrow$ QoS requirements vector
 $R \rightarrow$ resources matrix
 $r_j \rightarrow$ column of R

The Intelligence embedded in the SM

Workload Partitioning → Minimization Problem

$$\min_{\lambda} J(\lambda) = E|e_k|^2 = E|d - R\lambda_k|^2$$

subject to $\lambda(j) \geq 0, j \in [1, N], \mathbf{1}^T \lambda_k = \sum_{j=1}^N \lambda(j) = 1$

Adaptive Filter - Reweighted LMS

- Naturally fit for real-time estimation:** no need for previous training
- Able to **track variations** in the HPC resources (matrix R) and in the QoS requirements (vector d)
- Easy to be scaled up** towards several applications (kernels) and computing devices

$$\lambda_{k,i} = P[\lambda_{k,i-1} + \mu D_{\lambda} R_i^T e_k(i)] + F$$

$$P = I - \frac{1}{N} \mathbf{1} \mathbf{1}^T \quad F = \frac{1}{N} \mathbf{1} \quad D_{\lambda} = \text{diag}\{\text{entries of } \lambda\}$$

Results

- Test: variation in the QoS requirements**
 d changes while R remains fix
- Adaptive Filter is able to provide a **new workload partitioning at runtime → Fast reaction**

