

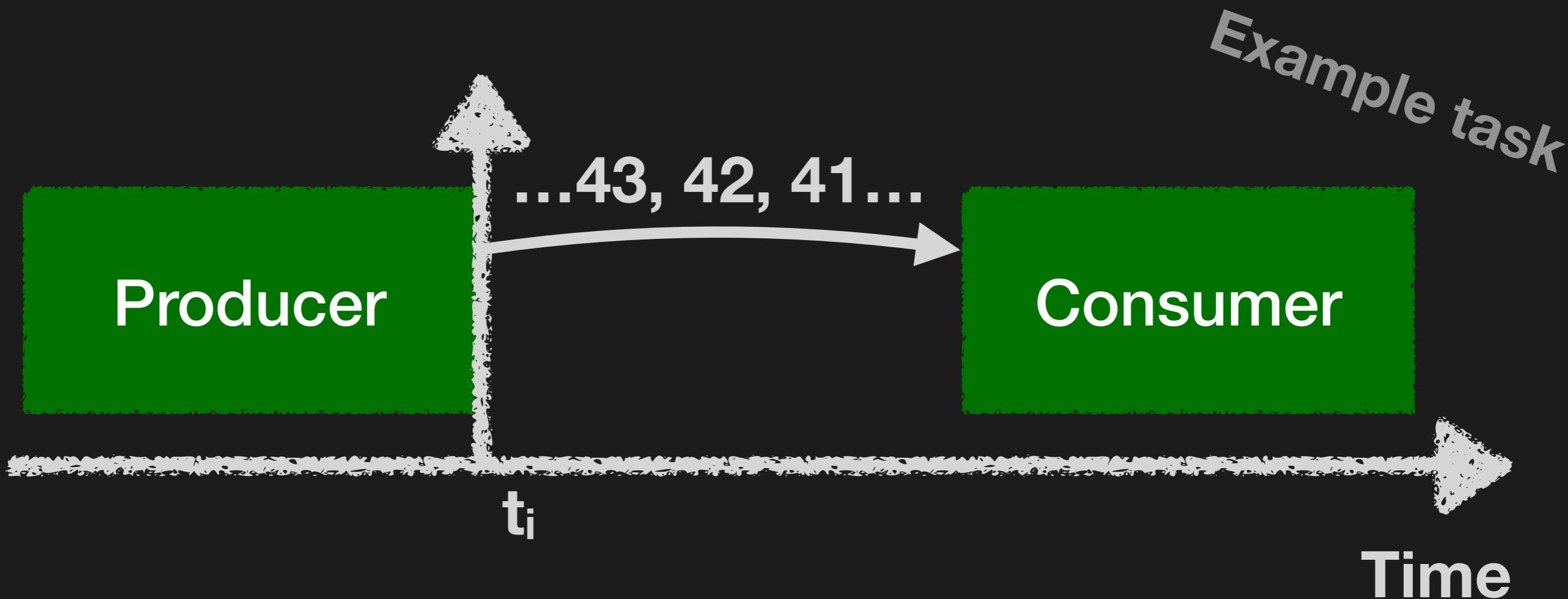
Enforcing Deadlines for Skeleton-based Parallel Programming

rtas2020.paulmetzger.info

**Programming parallel
real-time systems
is hard**

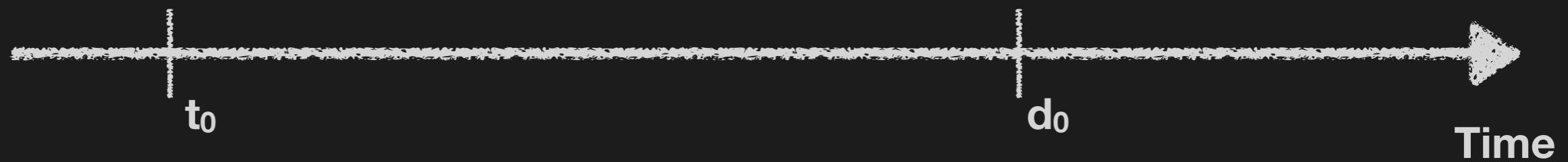
Parallel Real-Time Programming is Hard

A simple producer-consumer example



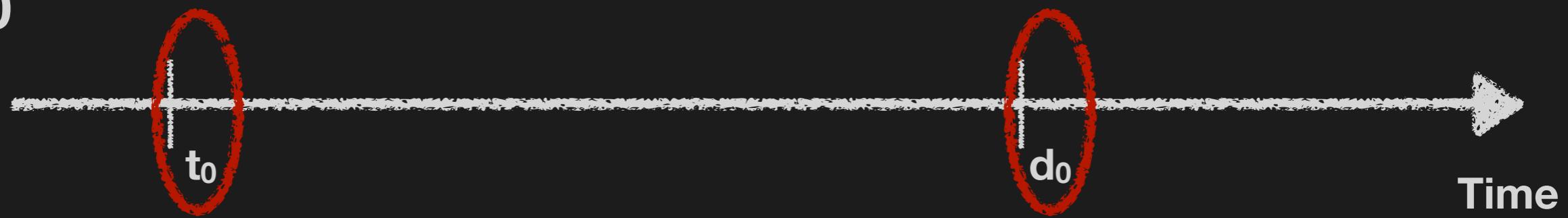
Parallel Real-Time Programming is Hard

Core 0



Parallel Real-Time Programming is Hard

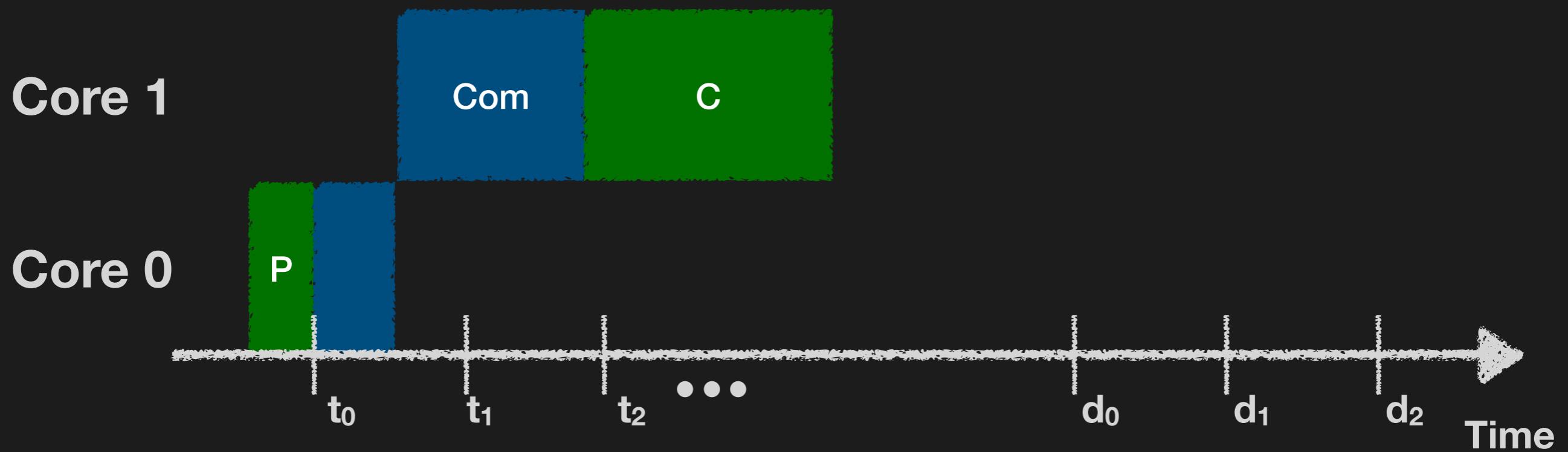
Core 0



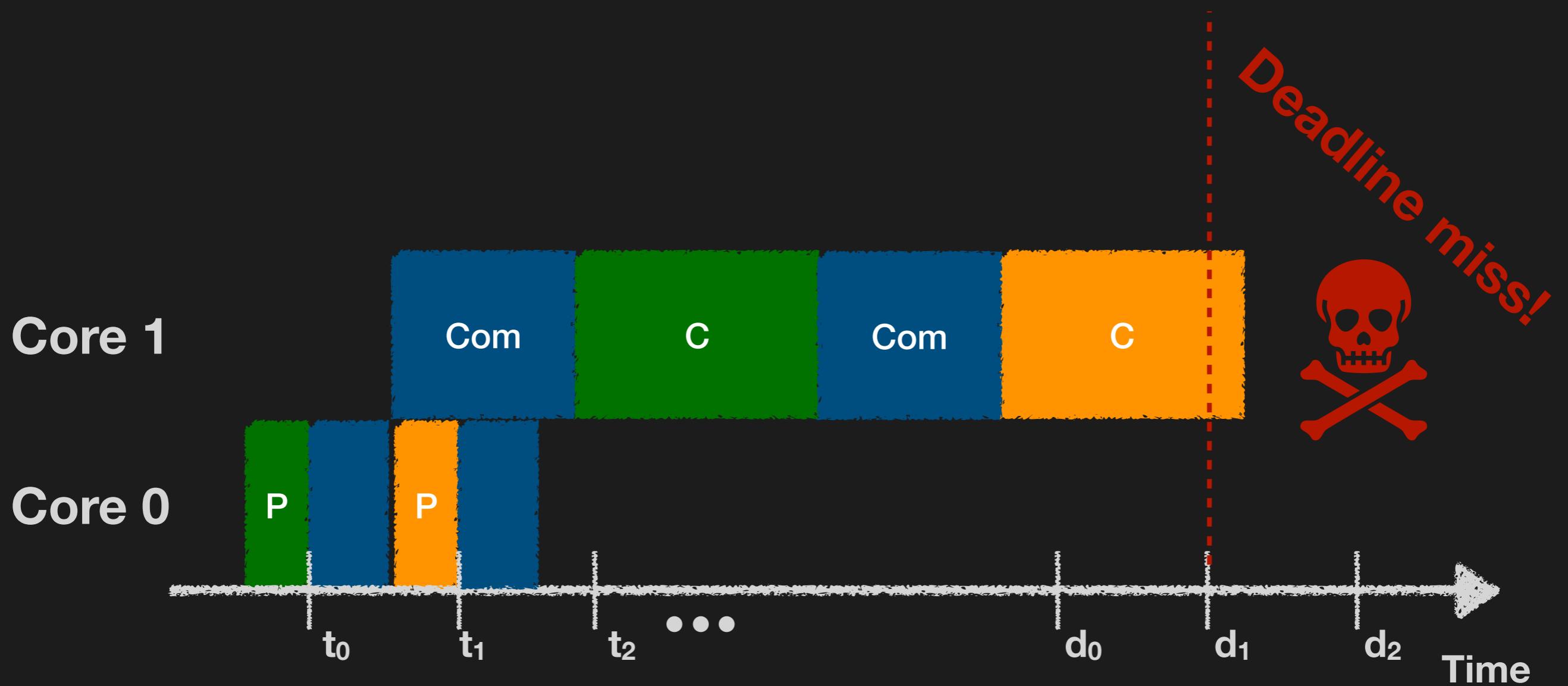
Parallel Real-Time Programming is Hard



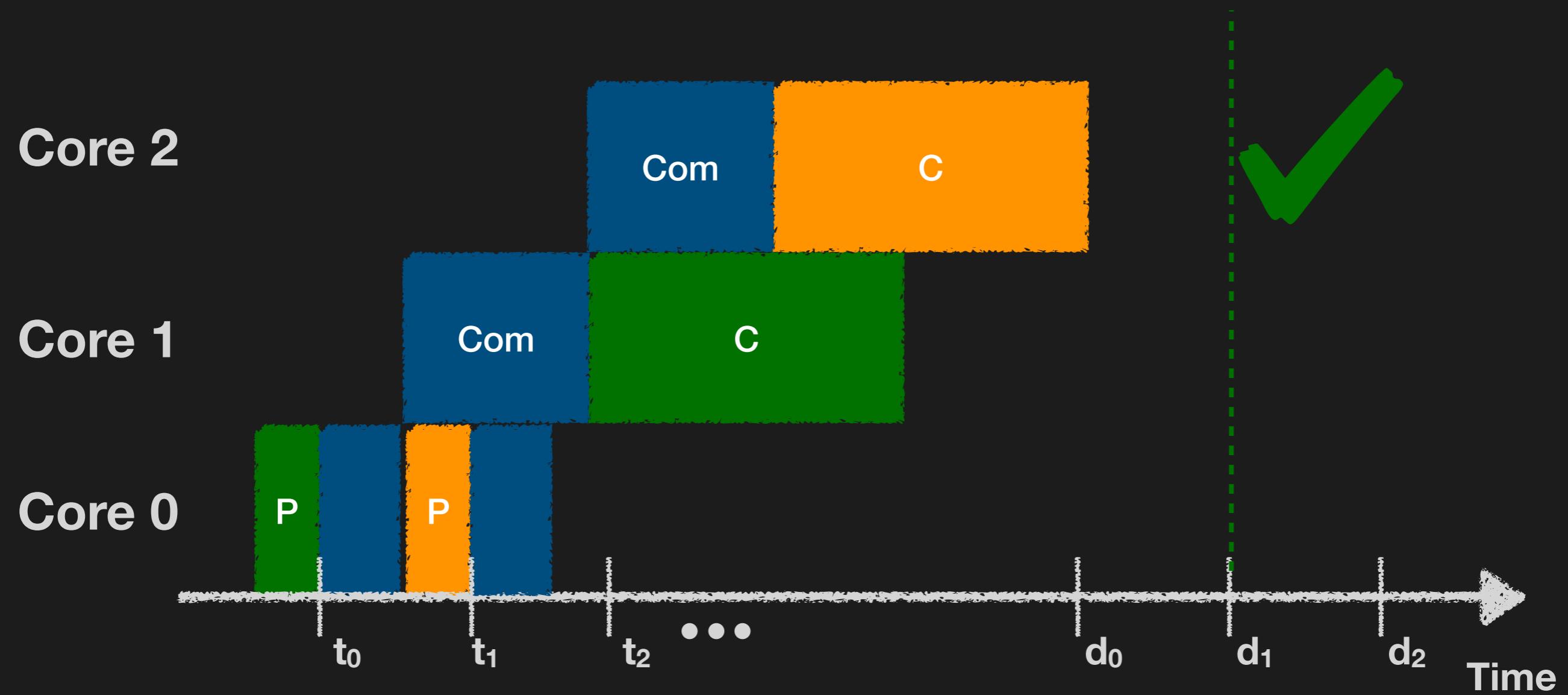
Parallel Real-Time Programming is Hard



Parallel Real-Time Programming is Hard



Parallel Real-Time Programming is Hard



Parallel Real-Time

We have to
determine the right
degree of parallelism

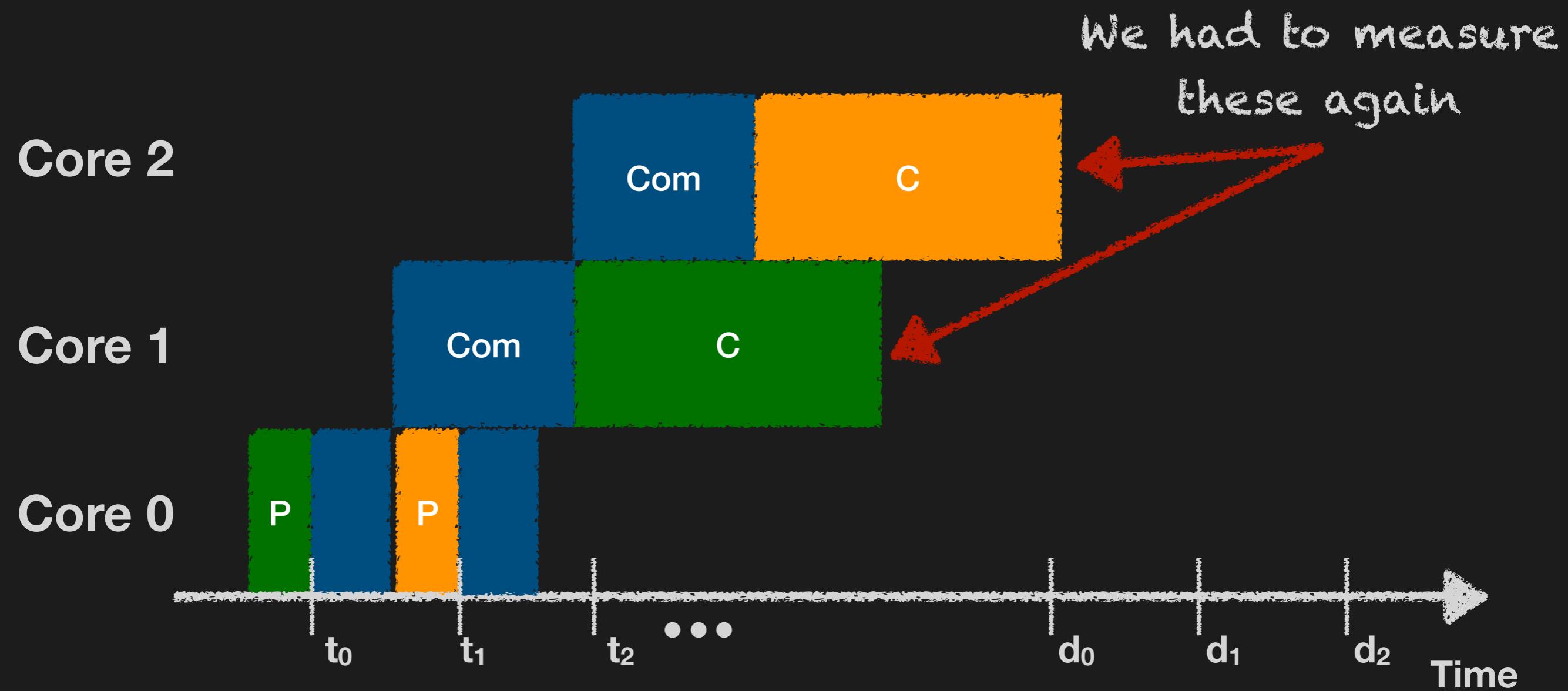
Core 2

Core 1

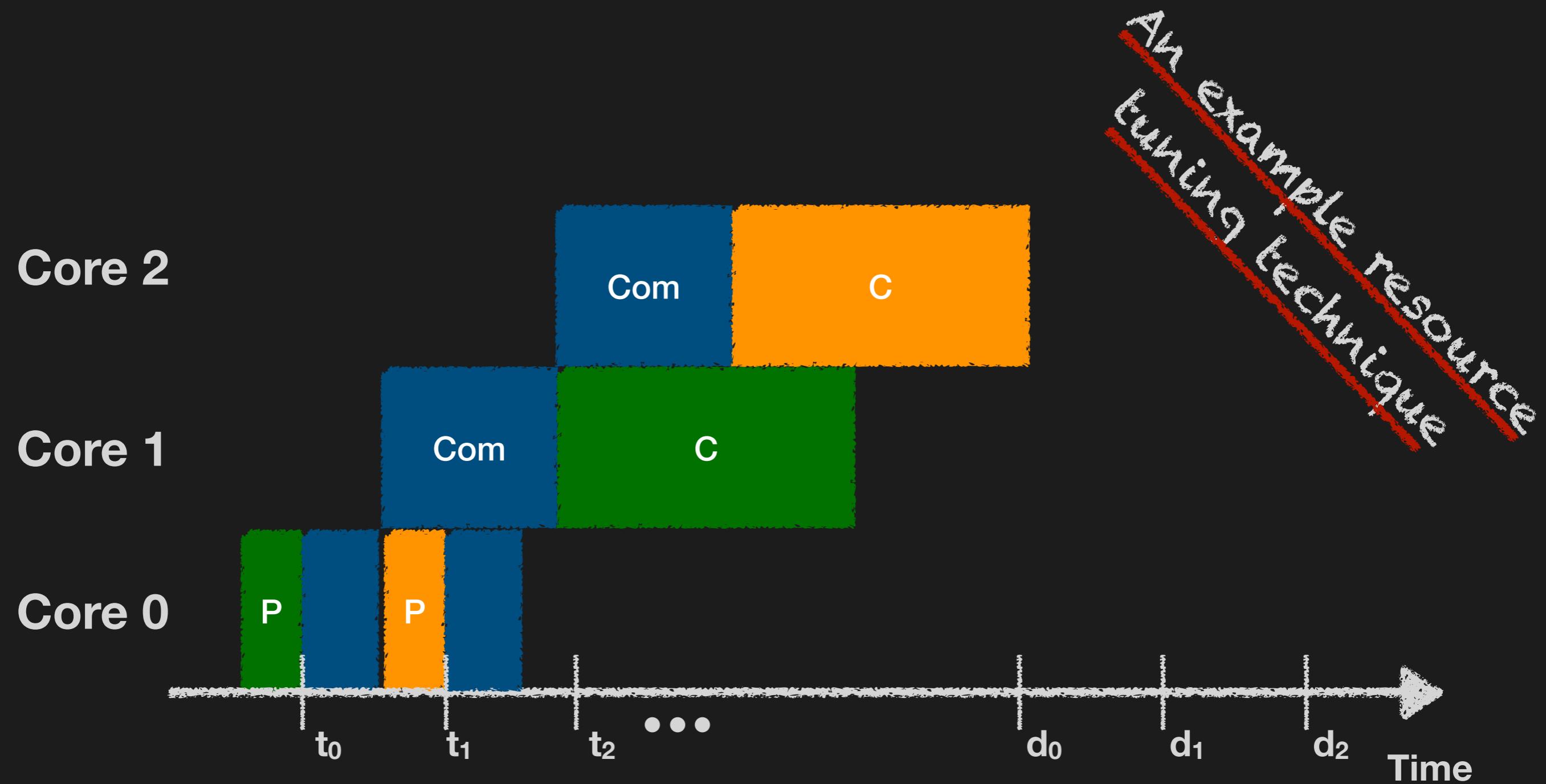
Core 0



Parallel Real-Time Programming is Hard

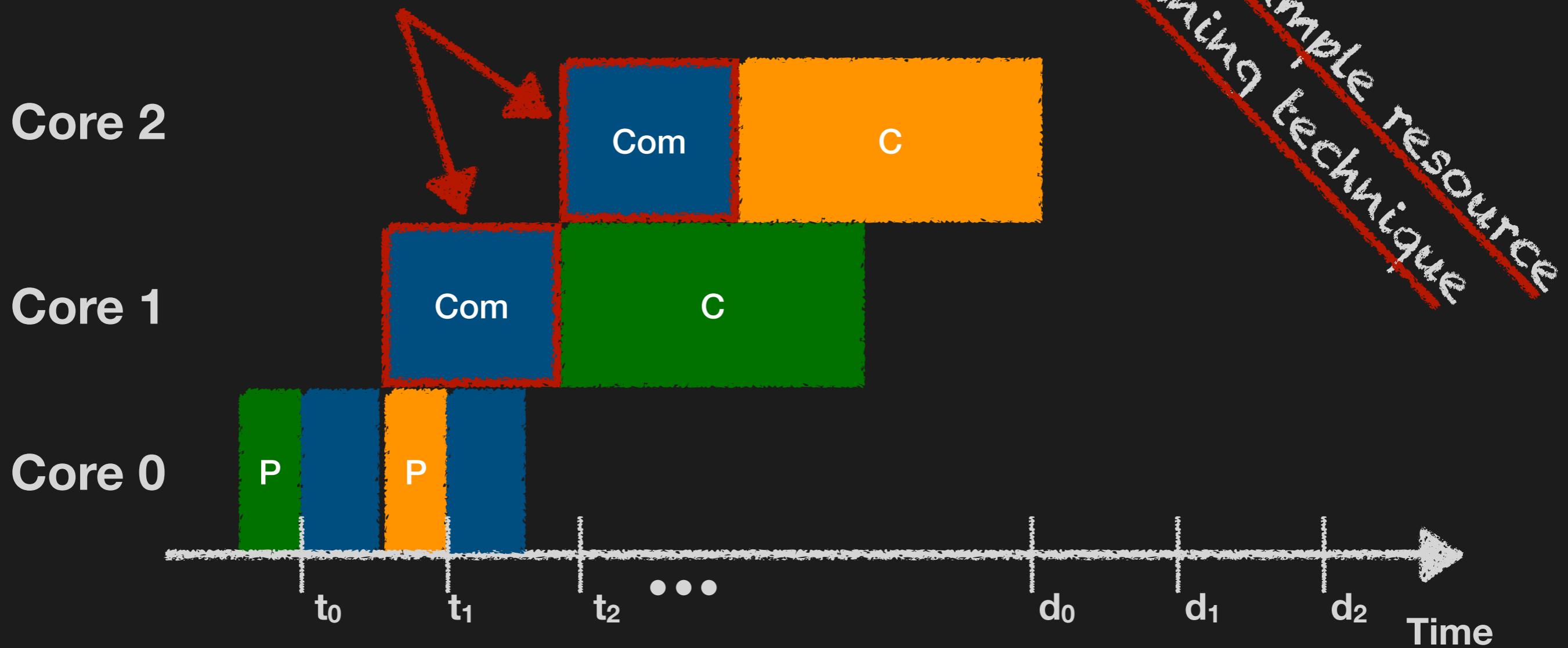


Parallel Real-Time Programming is Hard

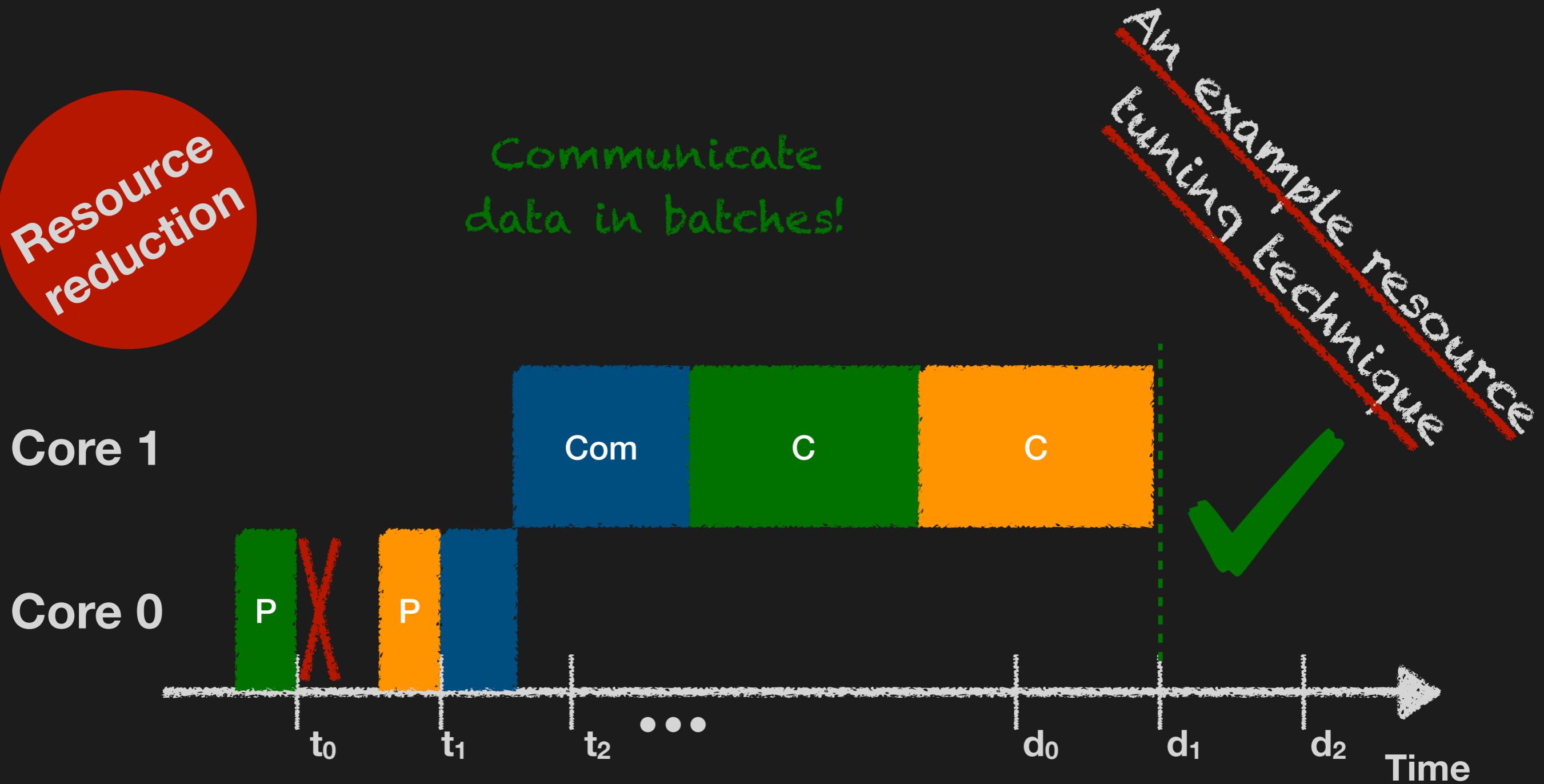


Parallel Real-Time Programming is Hard

We want to
reduce this

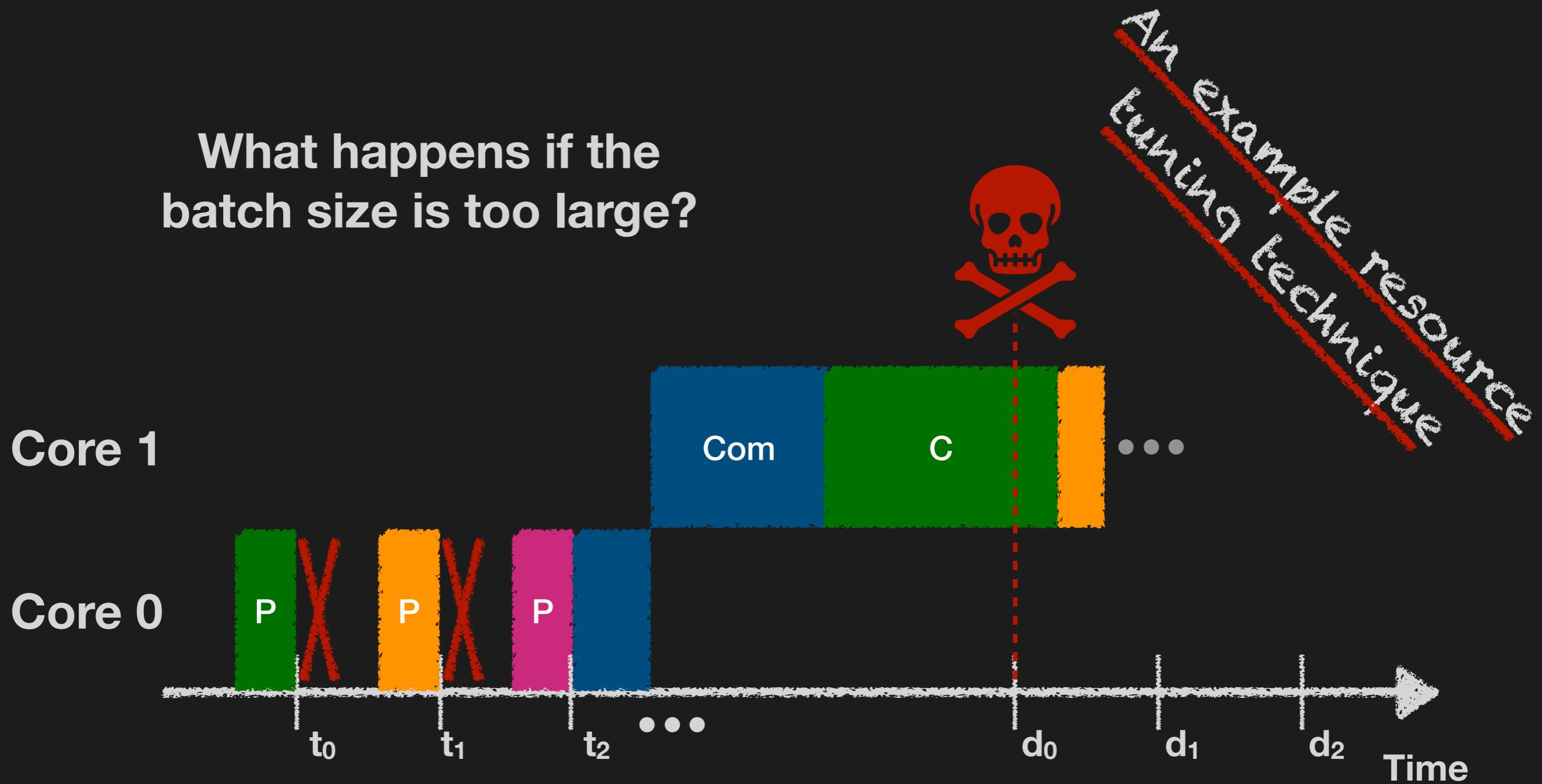


Parallel Real-Time Programming is Hard



Parallel Real-Time Programming is Hard

What happens if the
batch size is too large?

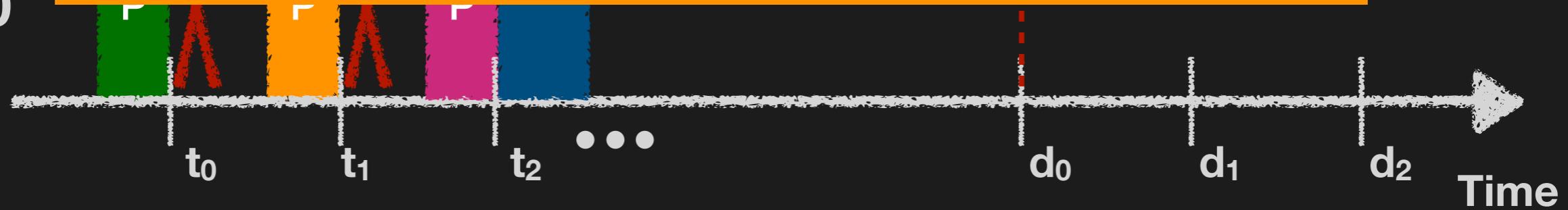


Parallel Real-Time

We have to
determine the best
batch size

Core 1

Core 0



~~idle resource
technique~~

What we want

- Automatically determine the degree of parallelism
- Automatically choose tuning parameters such as the batch size
- Hide complex parallel code from programmers

Our solution from another
planet  general purpose
parallel programming:

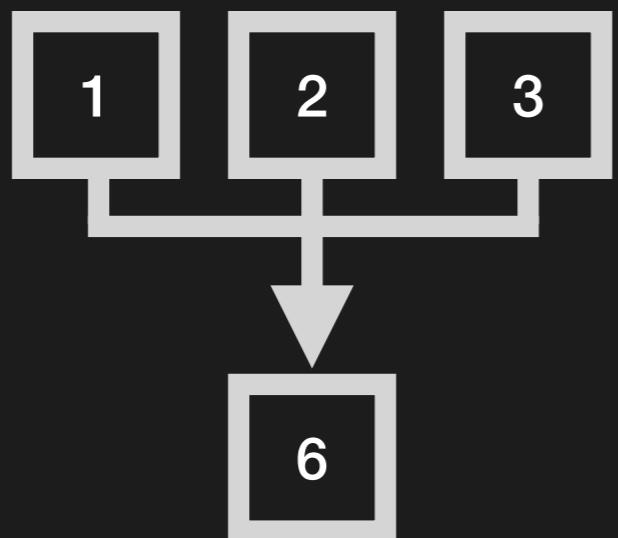
Skeletons

What are Skeletons?

Example 1

Reduction (+)

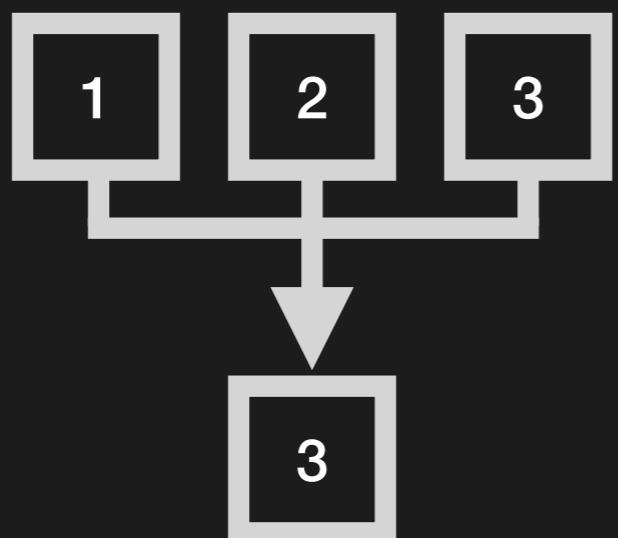
Input buffer



Output

Reduction (max)

Input buffer



Output

What are Skeletons?

Reduction (+)

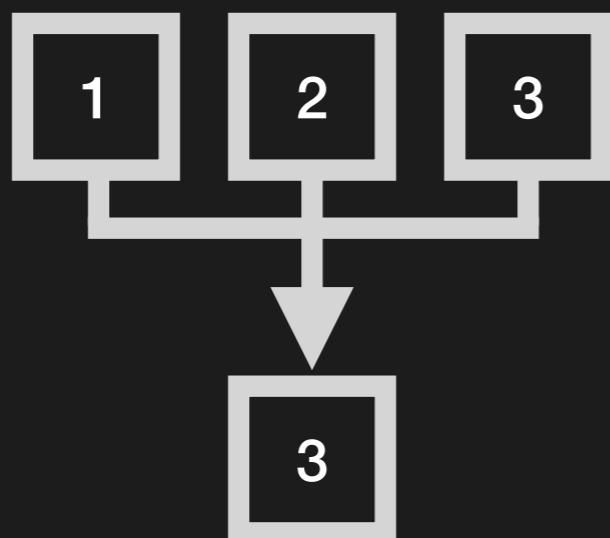
Input buffer



Output

Reduction (max)

Input buffer



Output

~~pthread_t t[THREAD_CNT];
for (...) pthread_create(...);
for (...) t[i].join();~~

~~int idx = ...
int offset = ...
for (...) {
 partial_result +=
 input_buffer[k];~~

~~for (...) result +=
 partial_result[j];~~

What are Skeletons?

```
reduction_skeleton(  
    input_buffer,  
    result,  
    +);
```

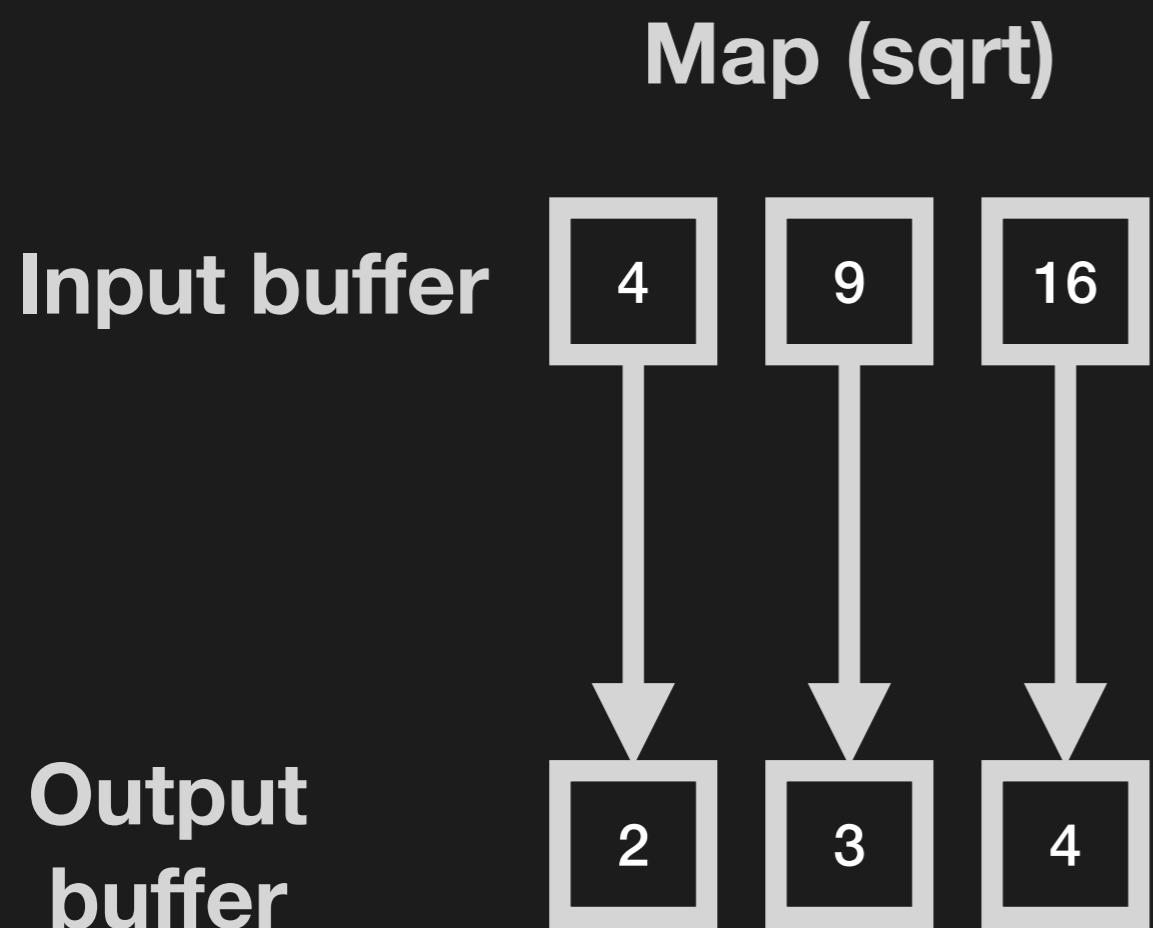
```
pthread_t t[THREAD_CNT];  
for (...) pthread_create(...);  
for (...) t[i].join();
```

```
int idx = ...  
int offset = ...  
for (...) {  
    partial_result +=  
    input_buffer[k];
```

```
for (...) result +=  
partial_result[j];
```

What are Skeletons?

Example 2

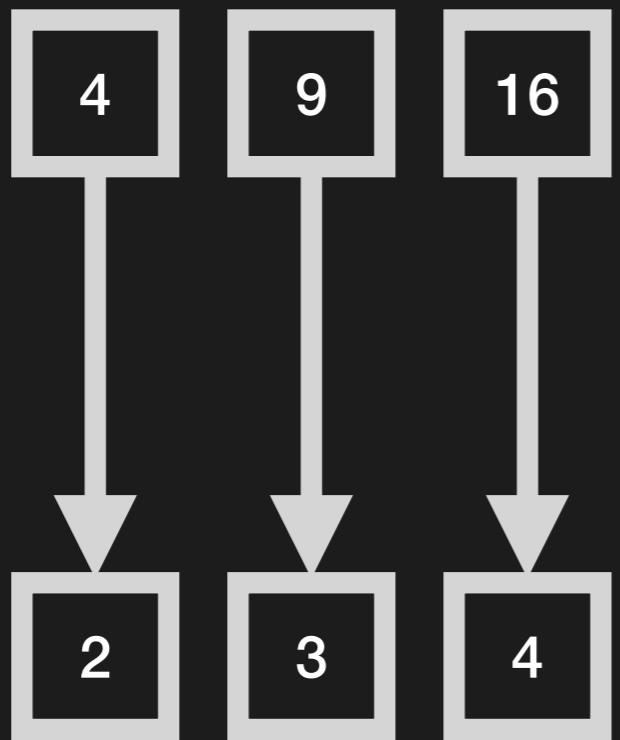


What are Skeletons?

Example 2

Map (sqrt)

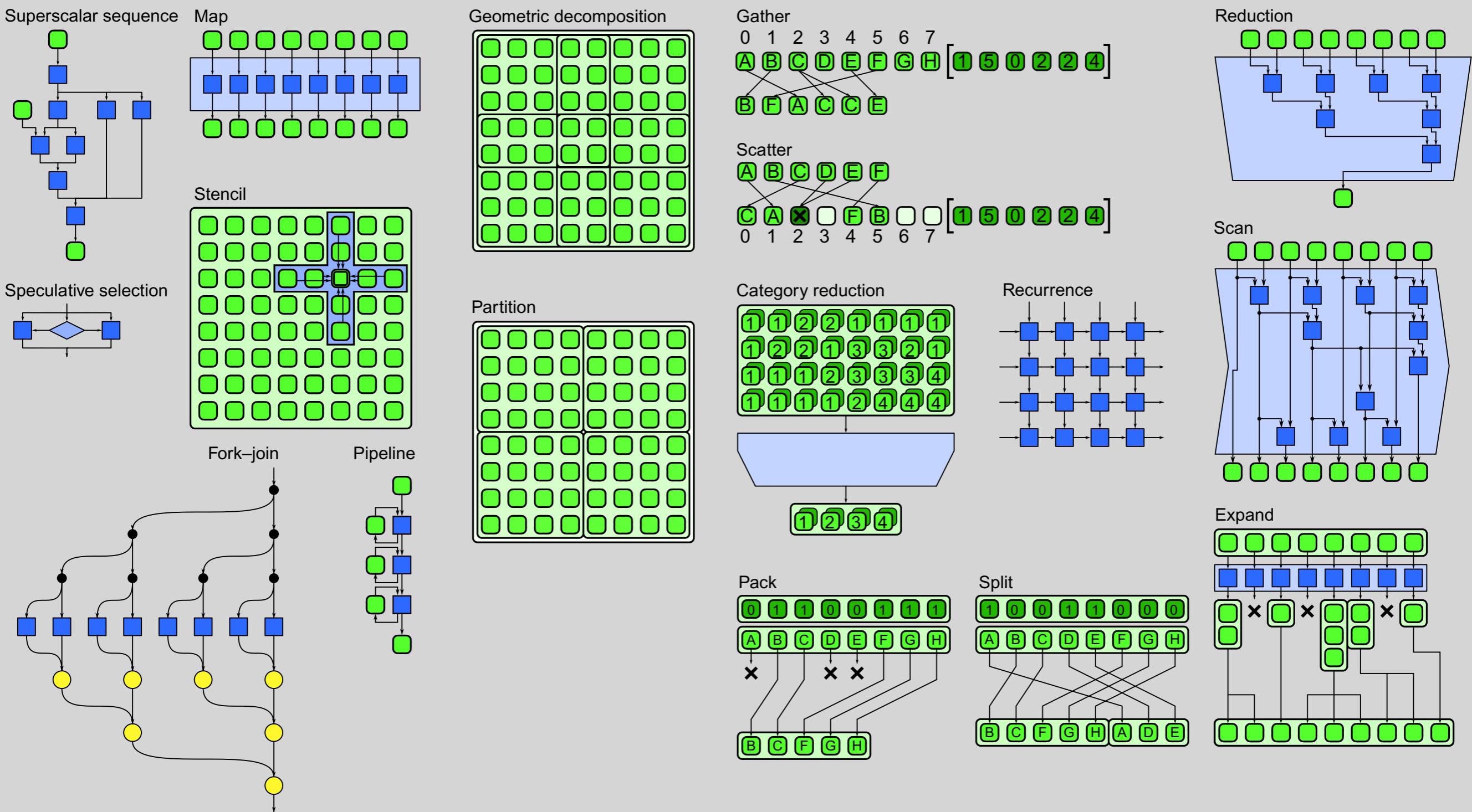
Input buffer



Output
buffer

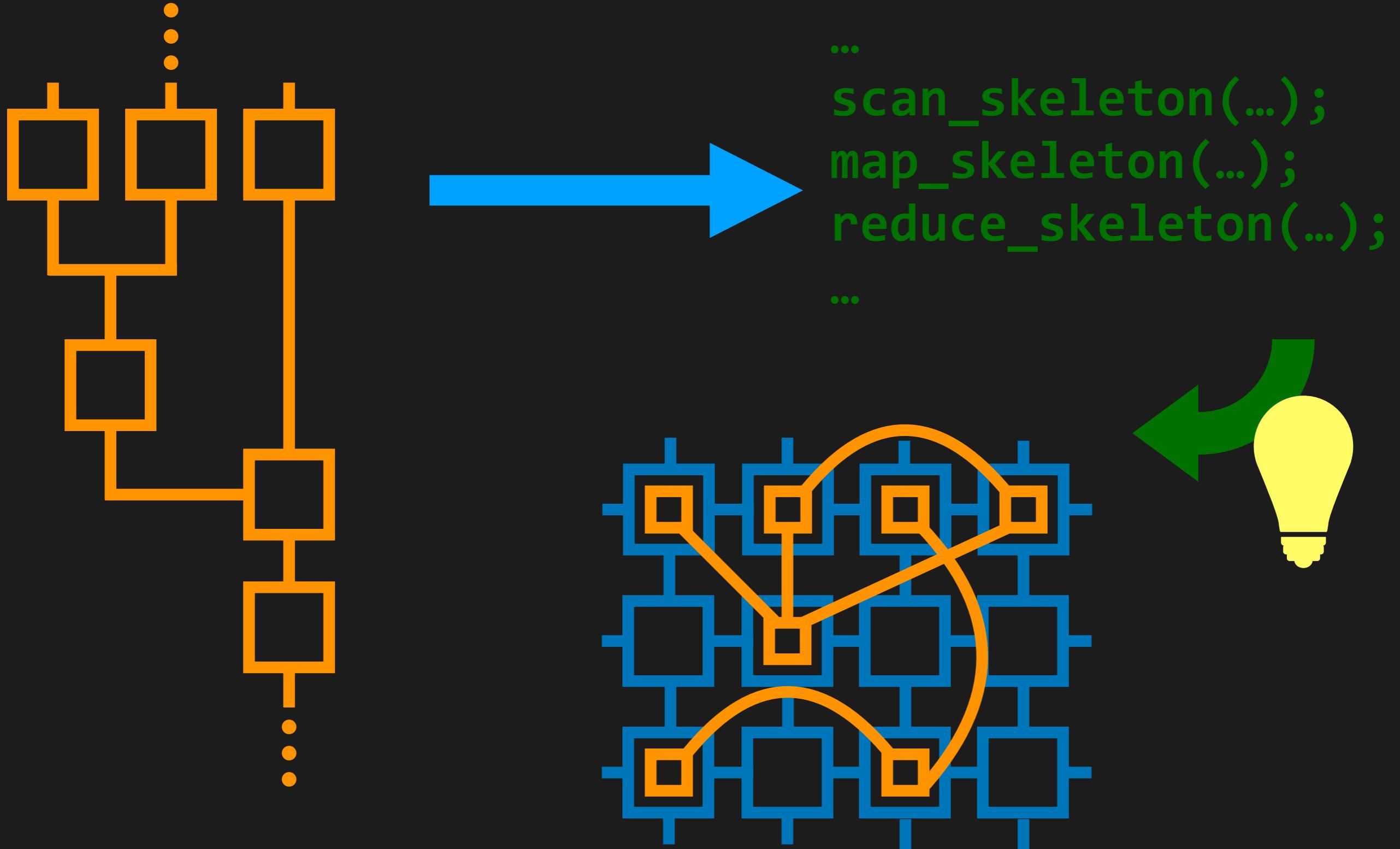
```
map_skeleton(  
    input_buffer,  
    output_buffer,  
    sqrt);
```

What are Skeletons?



from Structured Parallel Programming: Patterns for Efficient Computation. McCool et al., 2012

Skeleton Tuning



A case study with the producer- consumer skeleton

(We call this skeleton 'job farm' in the paper.)

Producer-Consumer Case Study



```
producer_consumer_skeleton(  
    producer_func, consumer_func,  
    consumer_wcet);
```

Producer-Consumer Case Study

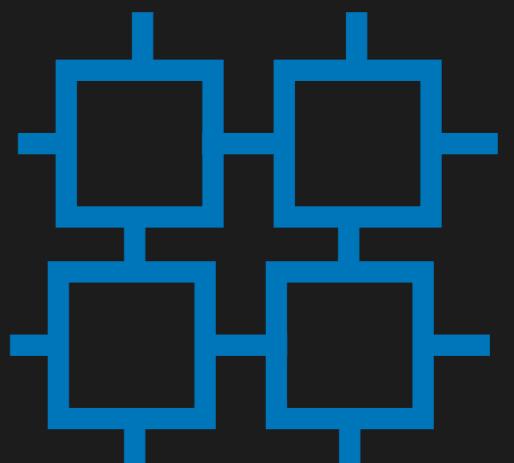
`producer_consumer_skeleton(...);`

**Ease of
programming**

```
int idx = ...  
int offset = ...  
for (...) {  
    partial_result +=  
        input_buffer[k];  
}  
  
...  
pthread_t t[THREAD_CNT];  
for (...) pthread_create(...);  
for (...) t[i].join();  
...
```



**Degree of
Parallelism**



Producer-Consumer Case Study

producer_consumer_skeleton(...); ↩

Set without
programmer intervention

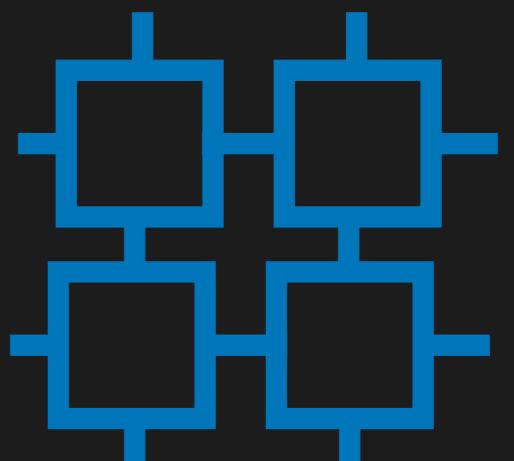
Ease of
programming

No implementation
specified

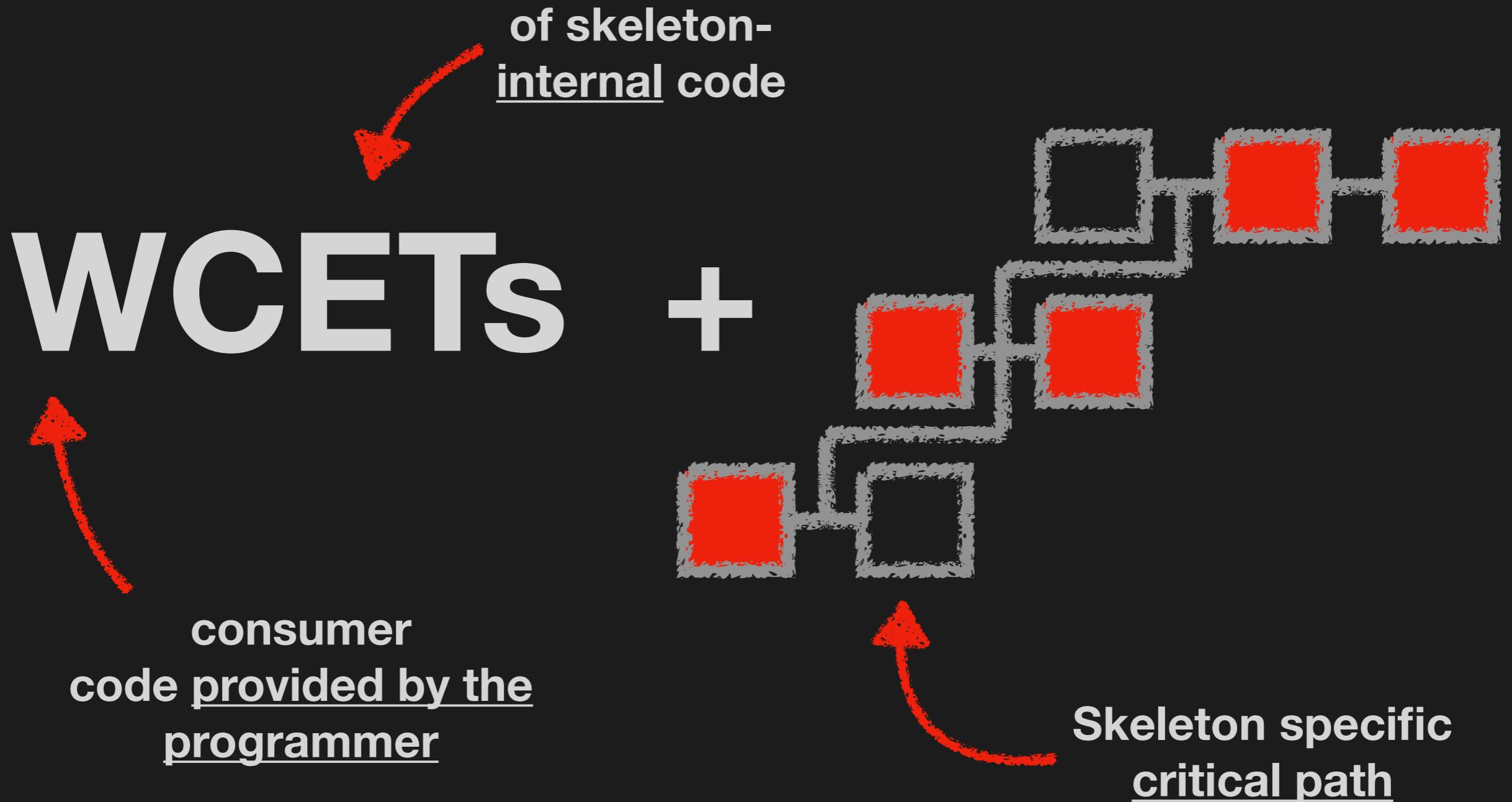
Batch size

Degree of
Parallelism

```
int idx = ...  
int offset = ...  
for (...) {  
    partial_result +=  
        input_buffer[k];  
}  
  
...  
pthread_t t[THREAD_CNT];  
for (...) pthread_create(...);  
for (...) t[i].join();  
...
```



How do we do it?

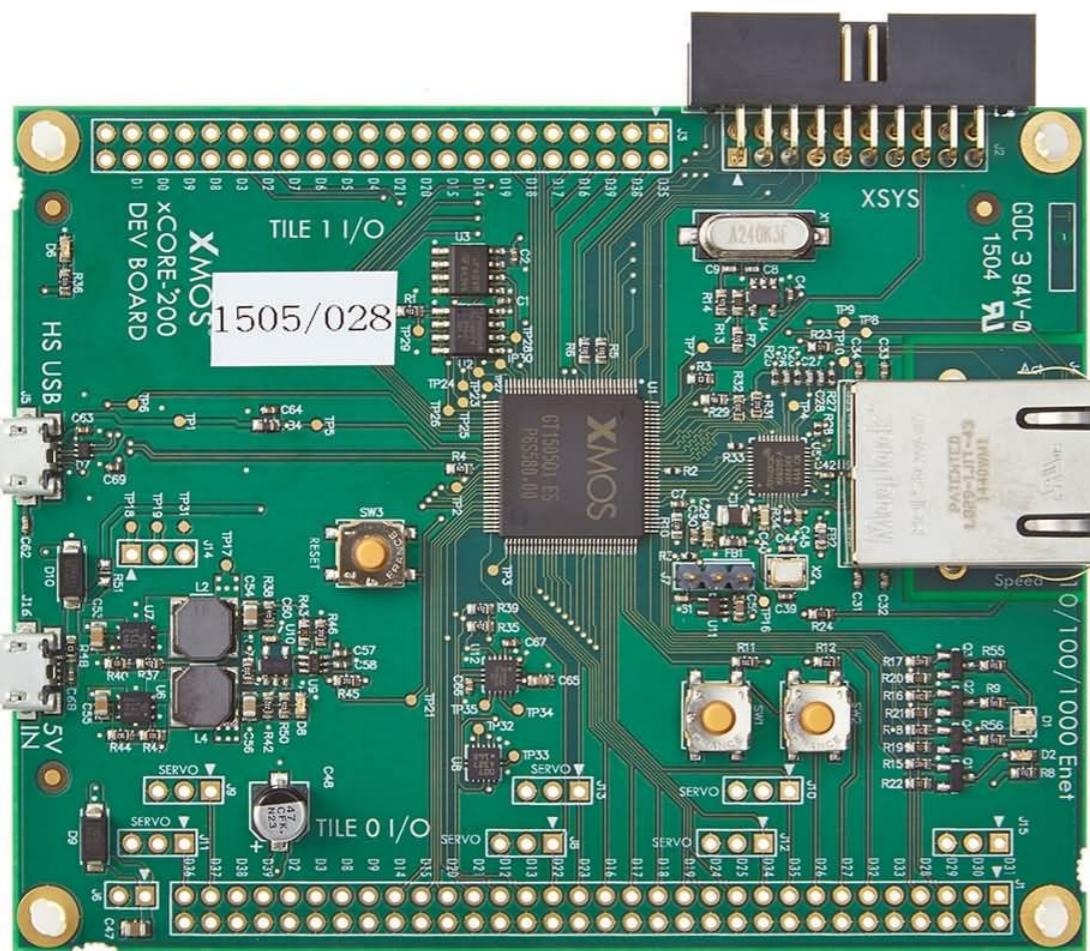


How do we do it?

Analytical
models for the { best DoP
best batch size

Results

XMOS xCore-200



Source: <https://www.xmos.com/wp-content/uploads/2019/04/explorerkit-200.jpg>

Real time:

- **Hard real time**
- **Static analysis for WCETs**
- **No data caches**
- **No OS**

Parallelism:

- **Fine grained multithreading**
- **Eight logical cores**
- **Cores share address space**

We find the best

degree of
parallelism in

89%

batch size in

72%

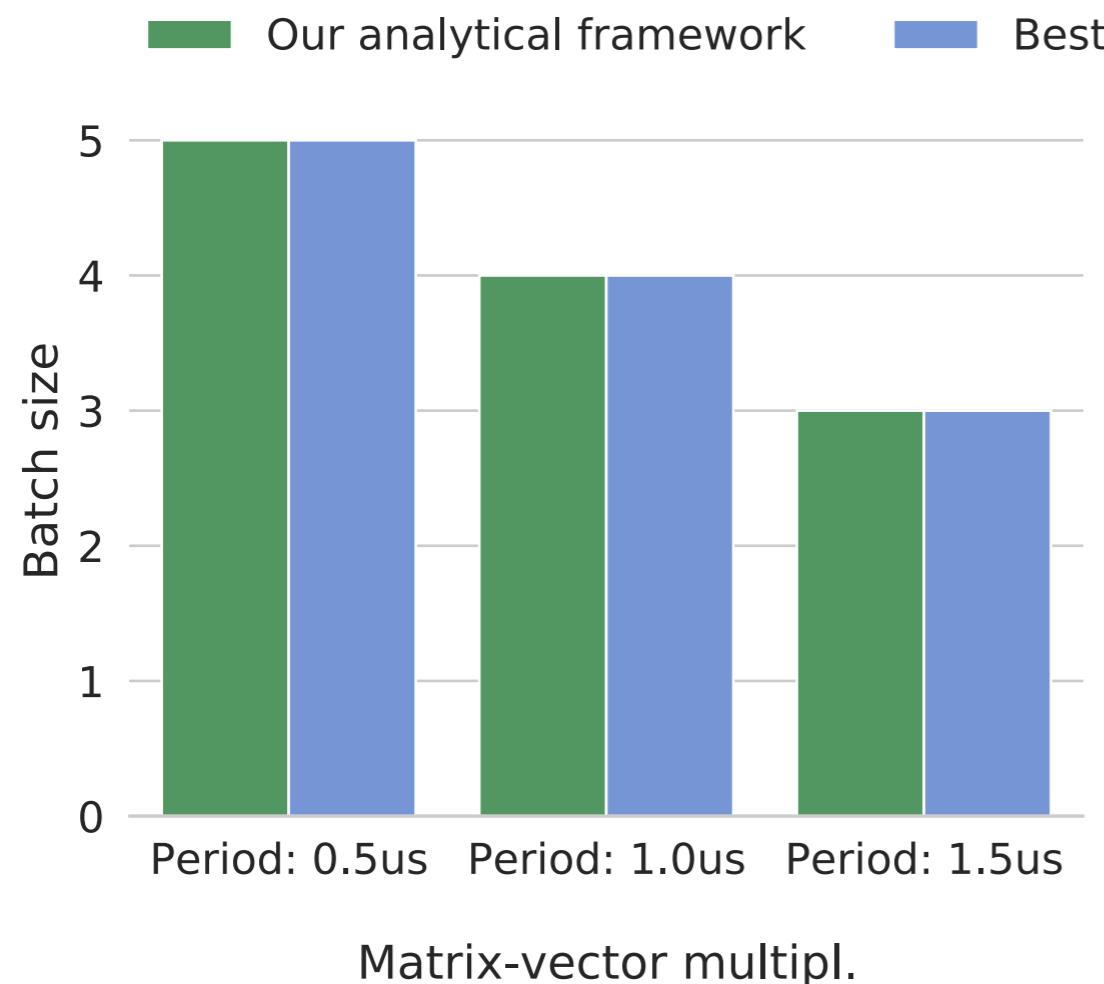
of all cases.

~~We use real
hardware~~

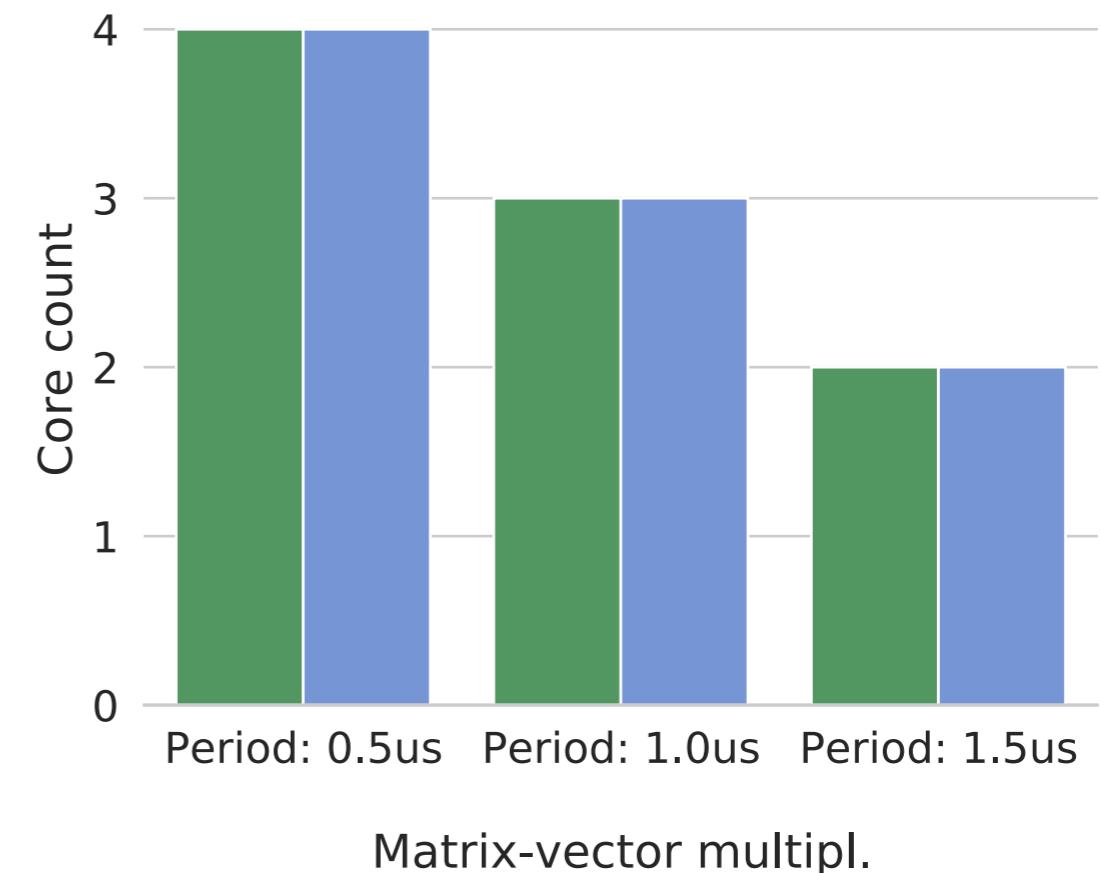
~~We never cause
deadline misses~~

Excerpt: Matrix Multipl.

Batch size choices

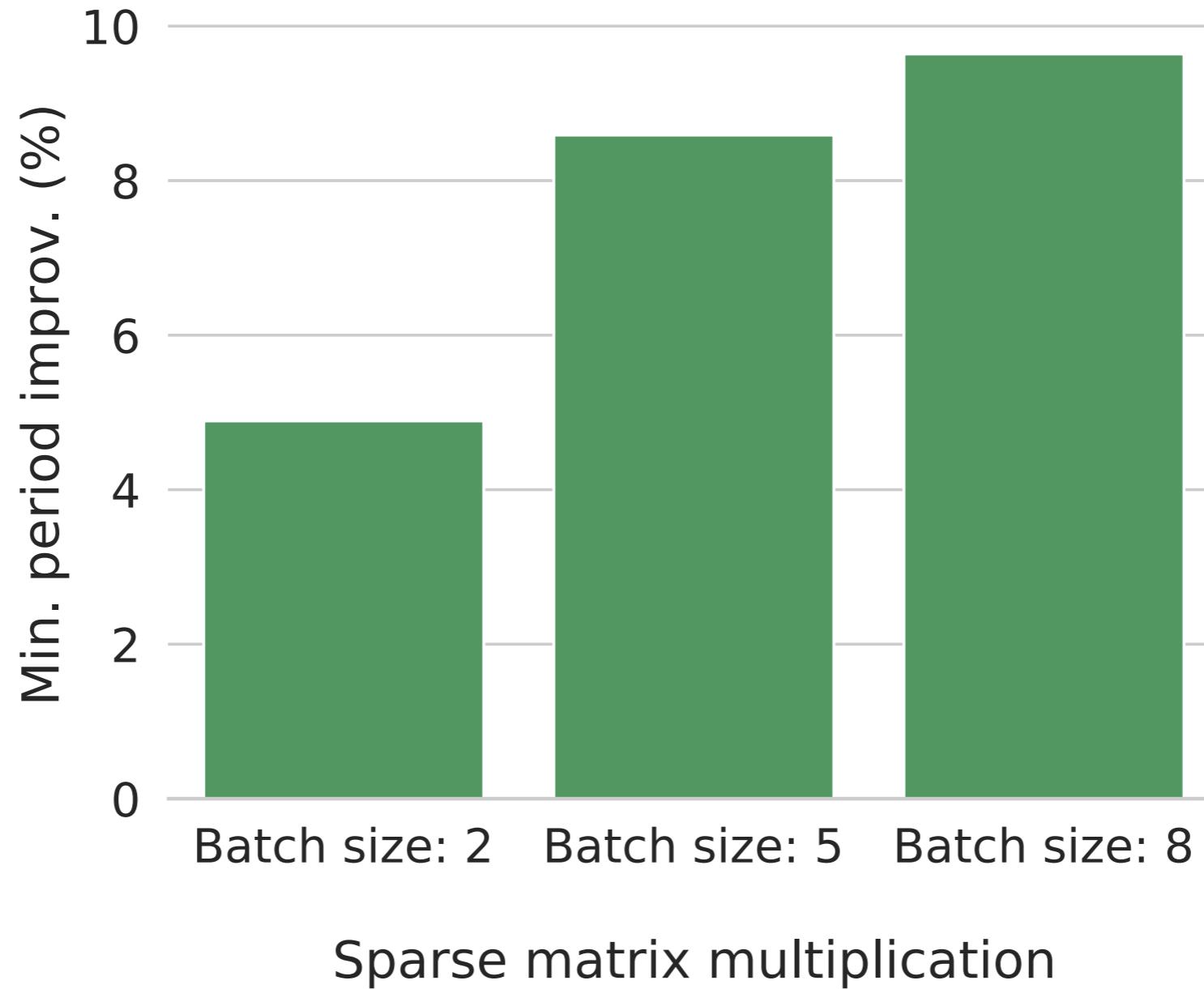


DoP choices



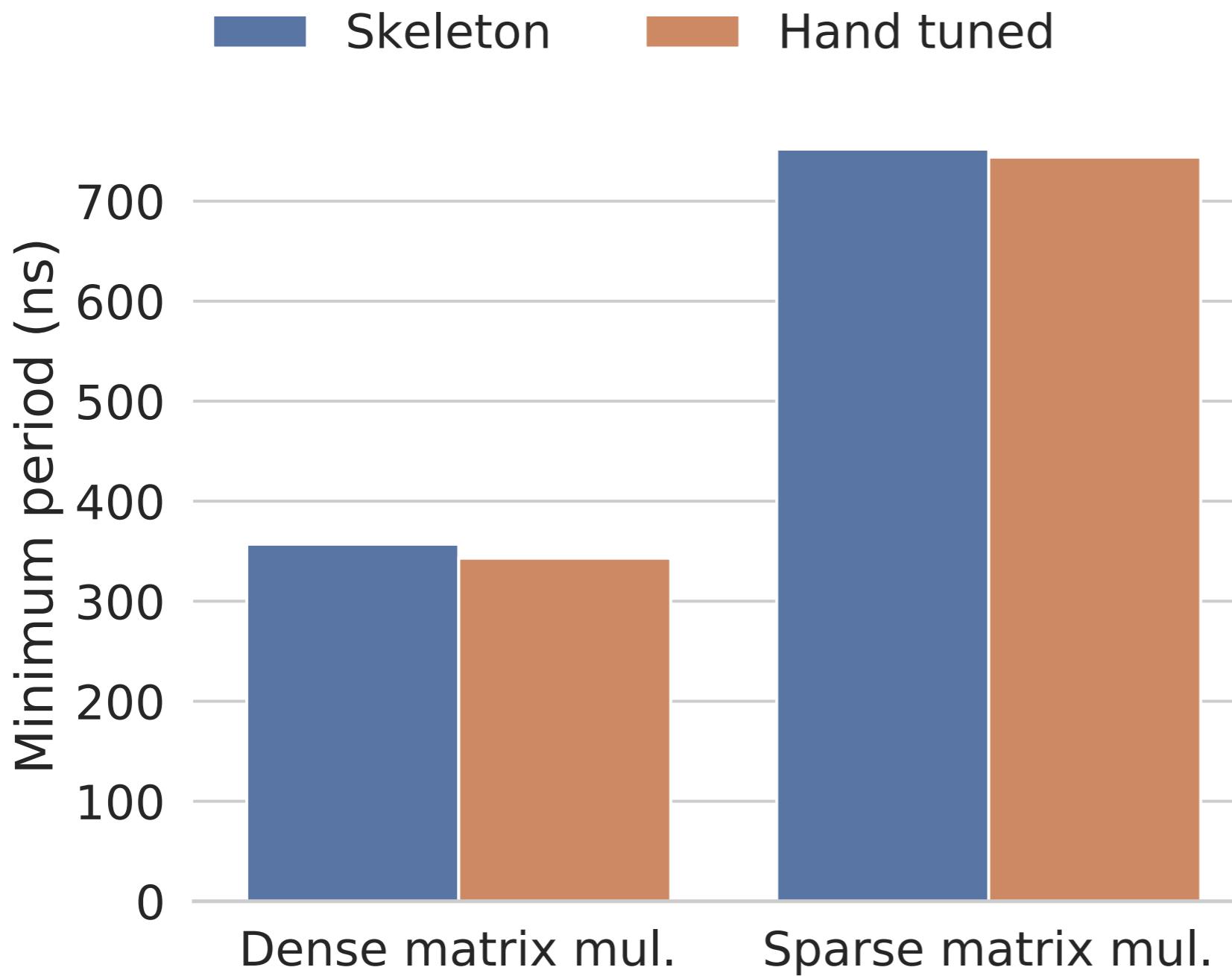
More benchmarks in the paper

Excerpt: Batching



Max. improvement with our benchmarks: 22%

Excerpt: Overheads



**Avg.
overheads:
3.48%**

What we want

- Automatically determine the degree of parallelism ✓
- Automatically choose tuning parameters such as the batch size ✓
- Hide complex parallel code from programmers ✓

Future work

More real-time skeletons

Being able to compose them

Enforcing Deadlines for Skeleton-based Parallel Programming

rtas2020.paulmetzger.info

Code is available online