



# ***Multiform Logical Time for MeMo-Codesign***

**Robert de Simone**

**Inria Sophia-Méditerranée & Université Côte d'Azur**

# INTRODUCTION

Today

**14th ACM-IEEE International Conference on Formal Methods and Models for (Embedded) System Design**

*Until 2013 was Formal Methods and Models for (Hardware-Software) Codesign*

**a) Hardware Design has traditional Models for Synthesis (up to a point)**

**b) System Engineering has its methodologies and diagrams**

**But Software Modeling and Formal Methods do not really match mainstream Software Engineering practices**

So what is relevant about Software Modeling for Embedded/CPS **system design**?

Or what could be ?

# OUTLINE

1. System-level, Model-driven, Platform-Based design
2. Around Hardware CAD and Synthesis
3. Around Software Engineering (*MoCs vs MoPs*)
4. A Clock Constraint Specification Language
5. Conclusion and discussion

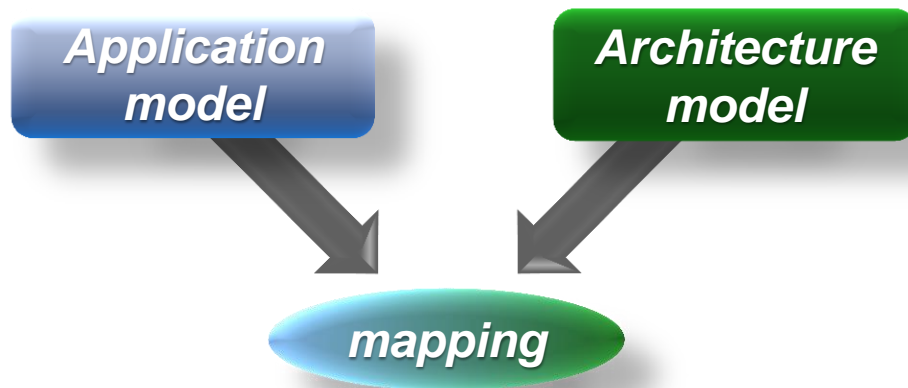
# 1

## **System-level, Model-driven, Platform-Based design**

**Multiviews, Cyber-Physical and Systems of Systems**

# Y-Chart design approach

Also branded as: Application / Architecture Adaptation (AAA)  
where adaptation means optimized mapping,  
and mapping covers allocation in space and scheduling in time

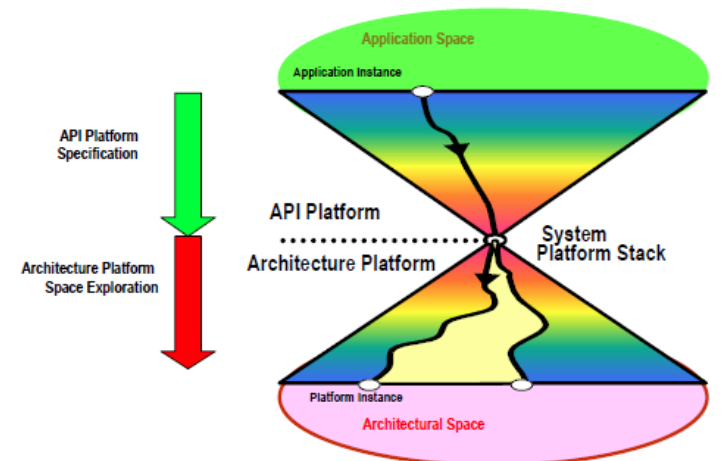
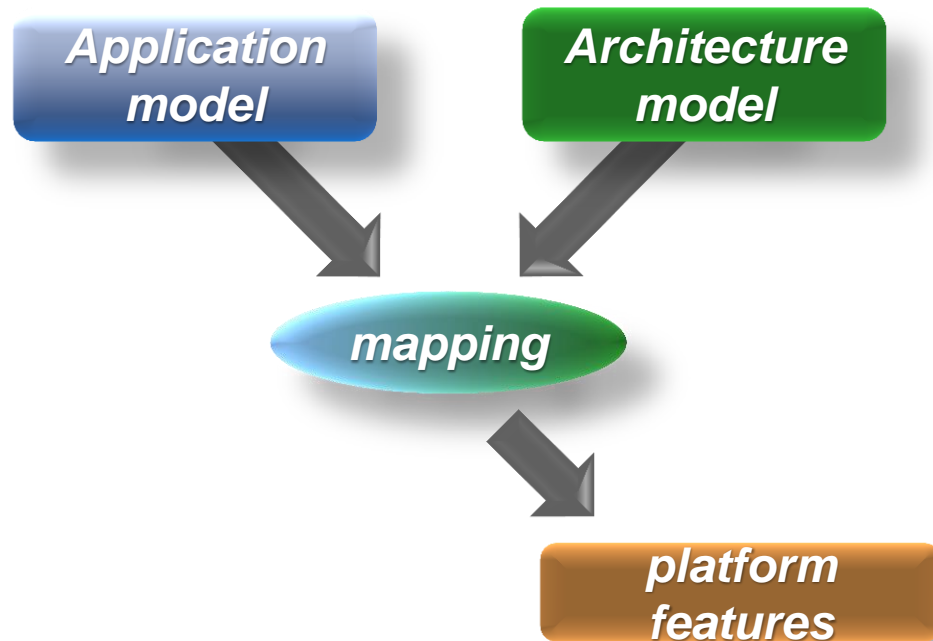


Actually used for design space exploration in distinct engineering purposes, at various modeling levels.

# Targeting Hardware: Platform-Based Design

Virtual Platforms and Virtual prototypes assembled by Hardware Architects.

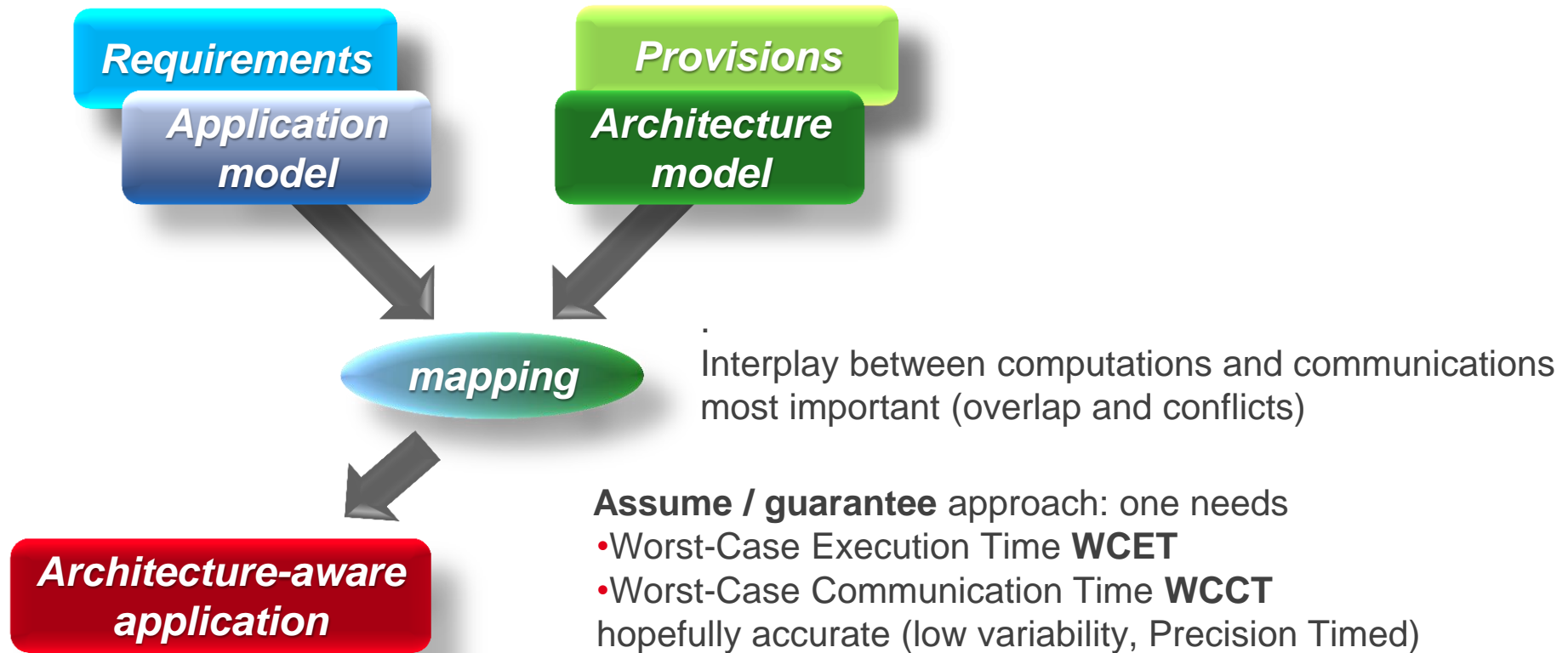
Application model is here *only* to provide typical use cases (booting Linux, one-day battery life) only to exercise the platform and accumulate non-functional information (typically by simulation in SystemC)



the version of Alberto Sangiovanni-Vincentelli  
(UC Berkeley + MARCO Giga-Scale Silicon Research Center)

# Application / Architecture Adaptation (AAA)

Here mapping corresponds to abstract compilation onto parallel heterogeneous architectures.

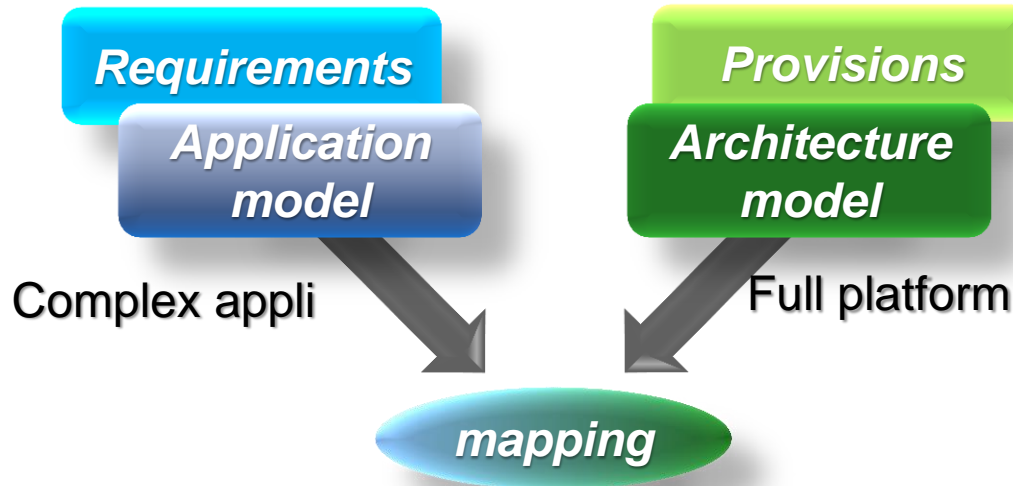


Typically used in **Real-Time Scheduling and Reactive programming**, to get logical timing guarantees

**Q: can clever mapping make a non time-predictable platform act more predictably ?**

# A two-storey approach

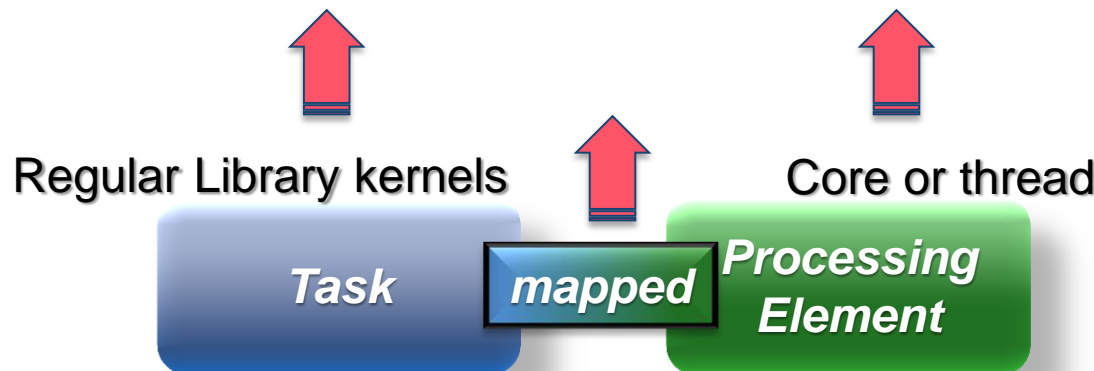
First Floor



The AAA approach cannot span the complete design

Predicable WCET and WCCT computed at a component level **(which?)**

Working at Ground-level more messy (less abstract)  
→ needs to be applied on regular subcases only.



**Level separation may fluctuate to adjust the approach**

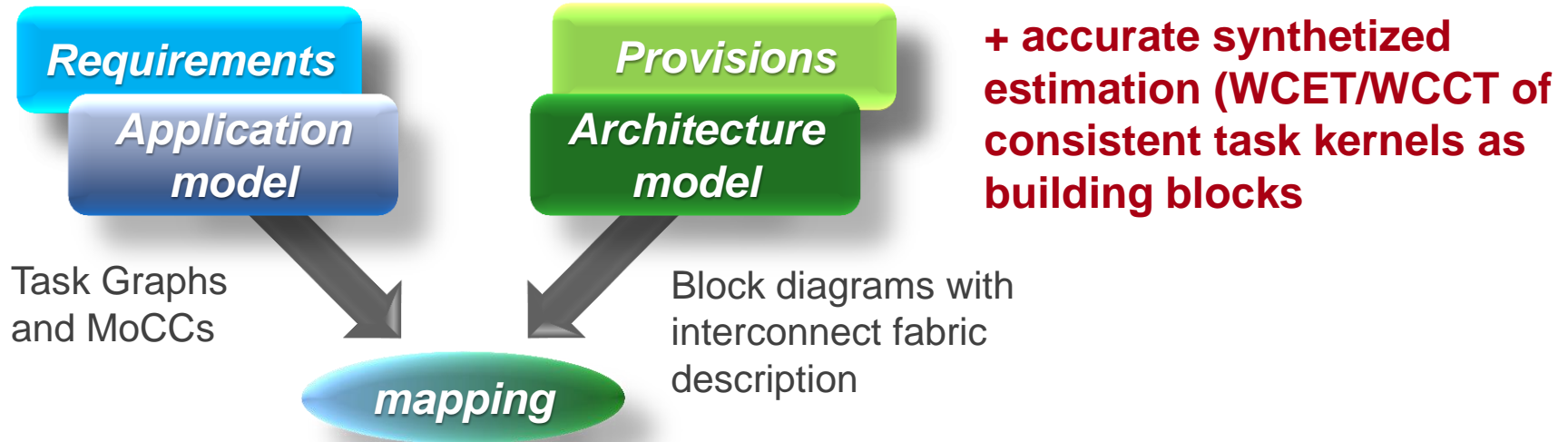
Ground Floor



# Adjust modeling to analysis (and vice-versa)

What we want:

Logical clock constraints + simple formulas between non-functional values



Automatic solver or abstract (regular) simulation for optimization

Before we can get this (in a realistic fashion?) we need more feedback  
How are similar things done elsewhere ?  
Can we borrow, compare, contradict ???

**Let's embark on a Tour.**

System Engineering

Formal Methods and Concurrency Theory (MoCCs)

(Domain Specific) Language design

Parallel compilation

Runtime execution and Optimization

Simulation and Worst-Case Execution Time

Hardware Abstraction Layer

# Example model-driven platform-based AAA environments

## Early system-level design stages (specification)

- SysML/MARTE (AADL)
- ARCADIA / CAPELLA
- Amalthea

## Hardware Virtual platforms (MoCs and SoCs)

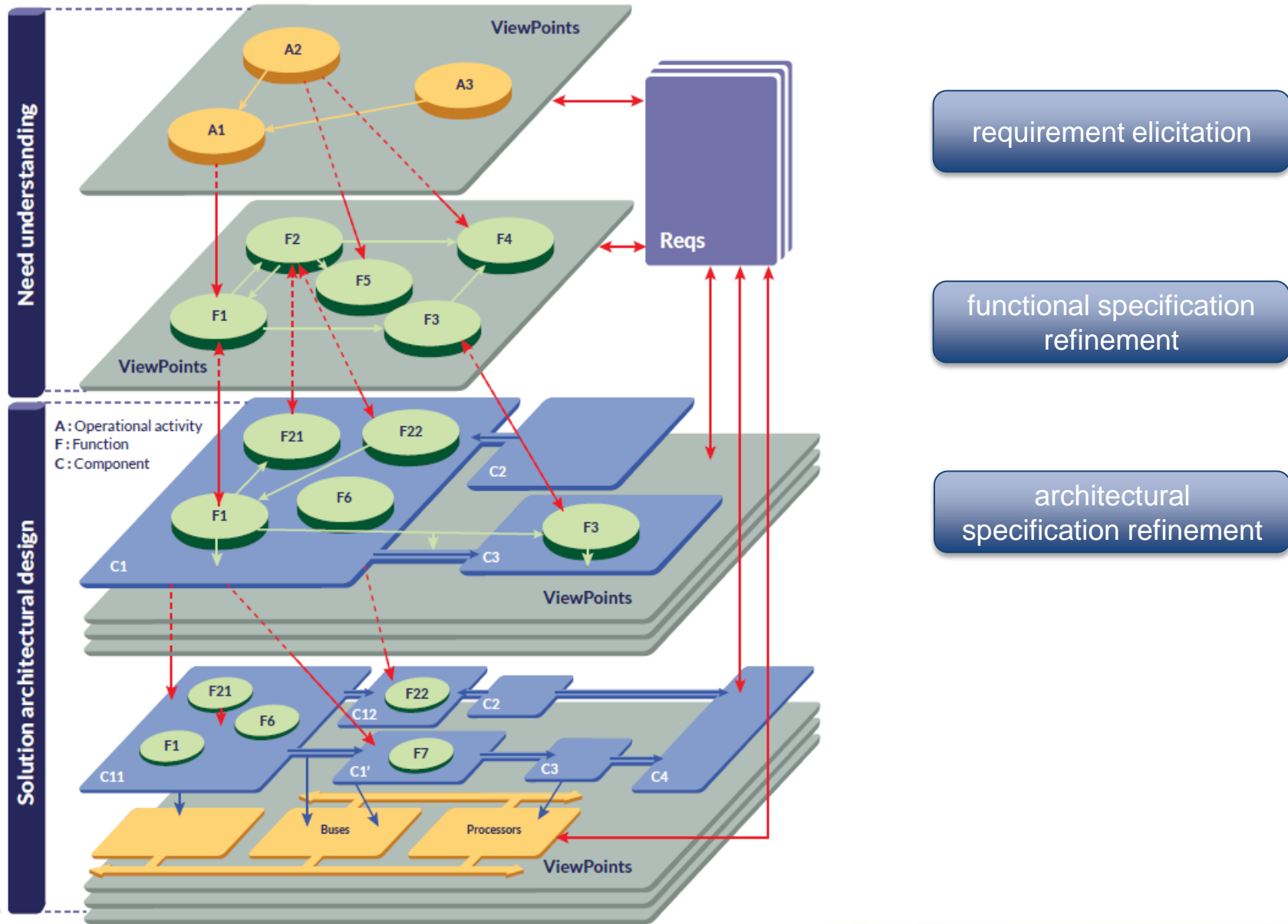
- Synopsys MCO Platform Architect
- *MetroPolis*
- *StreamIt/Raw (Tilera)*
- *SDF3/Aelite*
- *ForSyDE/Nostrum*
- *SigmaC/Kalray*

## Real-Time Scheduling

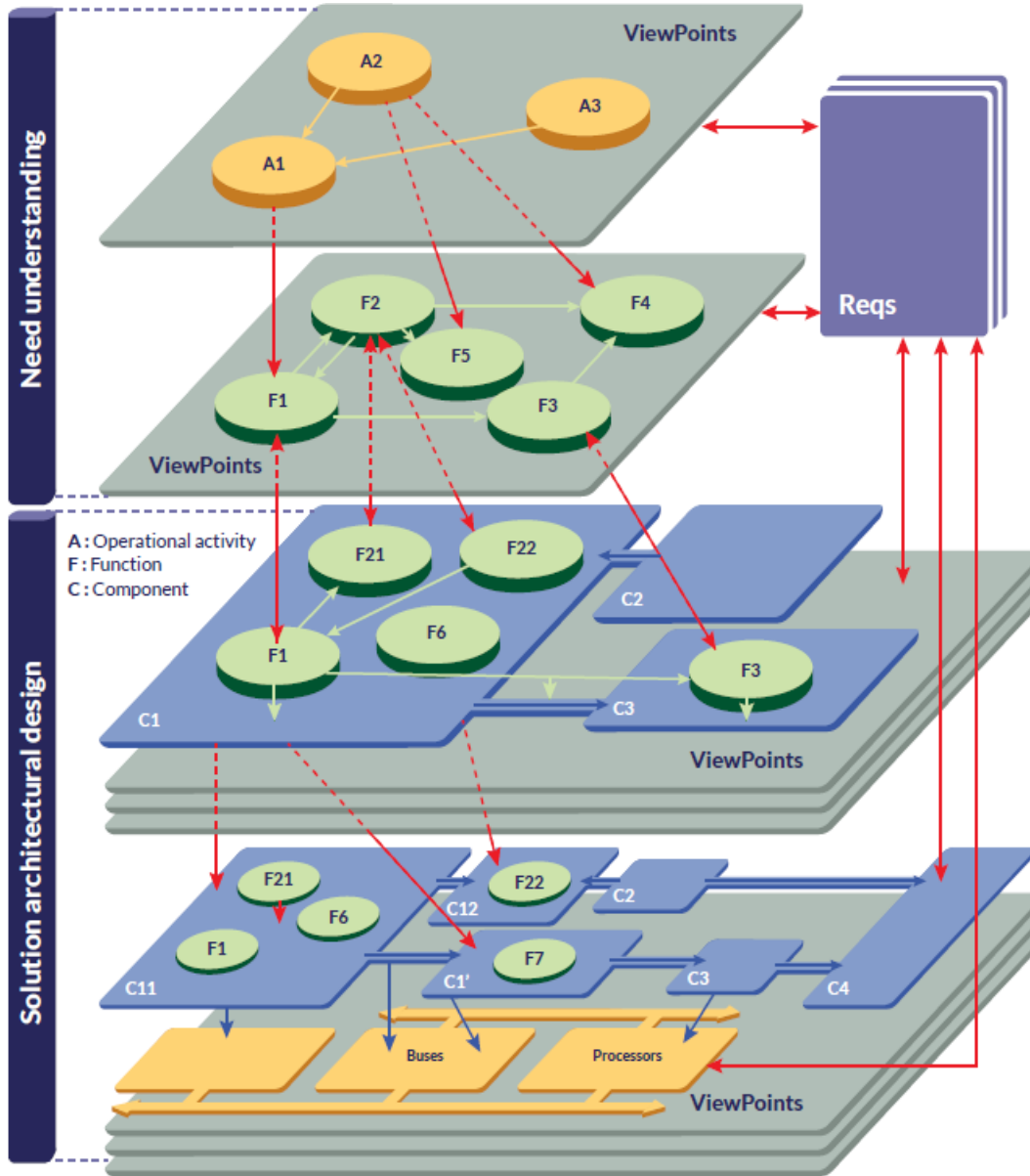
- AADL (SysML/MARTE)
- SynDEx

Much more often, the combination remains implicit...

# System Engineering Design Flow: Arcadia / Capella example



## System Engineering Design Flow: Arcadia / Capella example



Allocation made by user, on  
software and hardware architectur

Quality of allocation evaluated by computing **simple cost functions**:

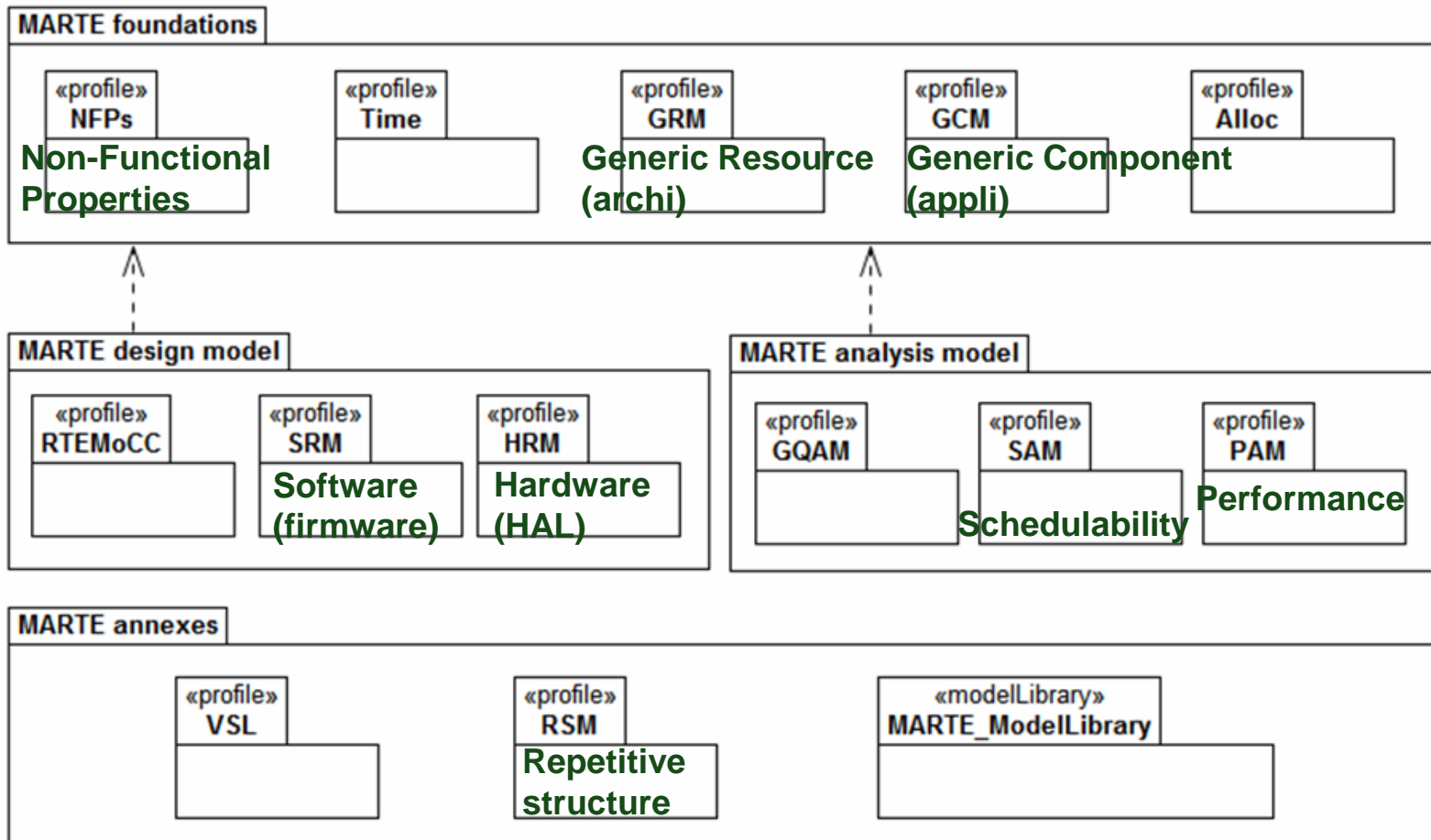
## Excel-like spreadsheets mostly

## formulas $\leftrightarrow$ constraints

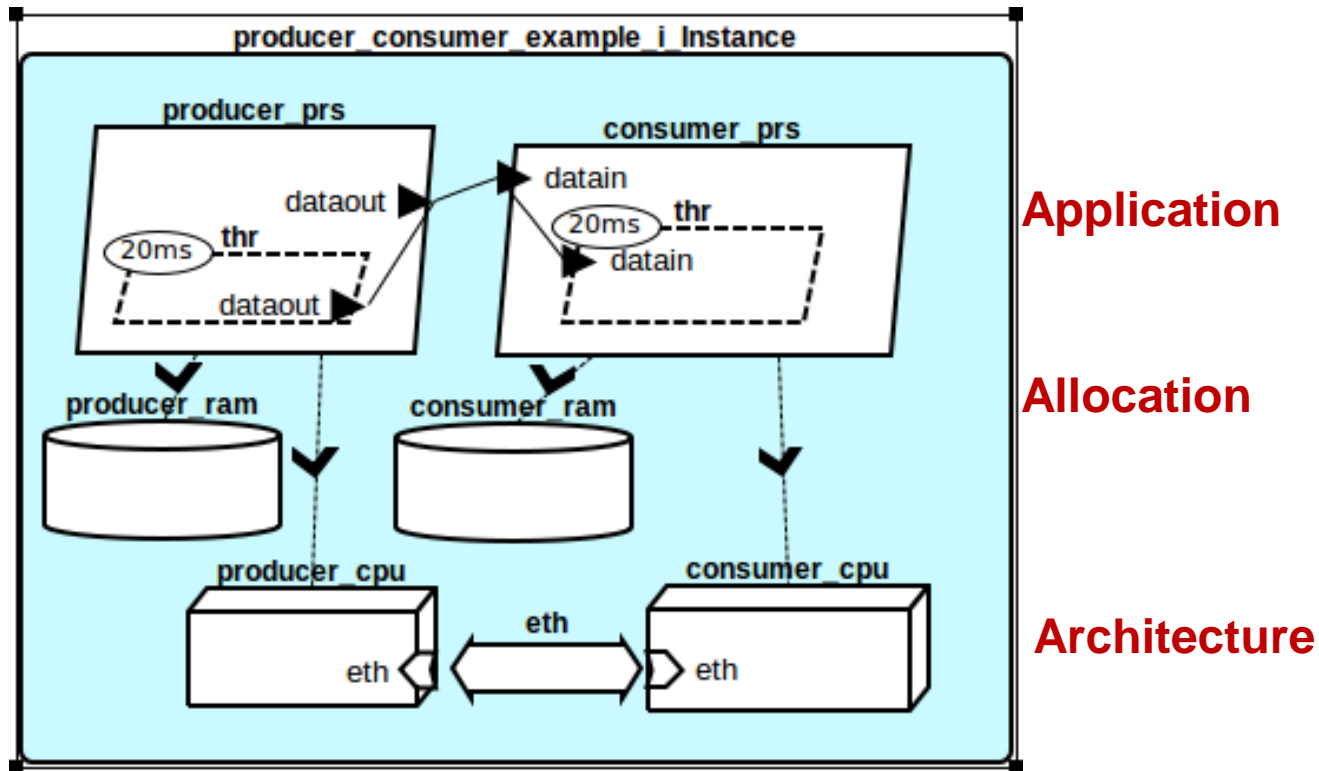
Example:

- (money) cost
- mass
- reliability (fault tolerance)
- security (mixed-criticality)

But what if dynamics involved ?  
Mode&State changes



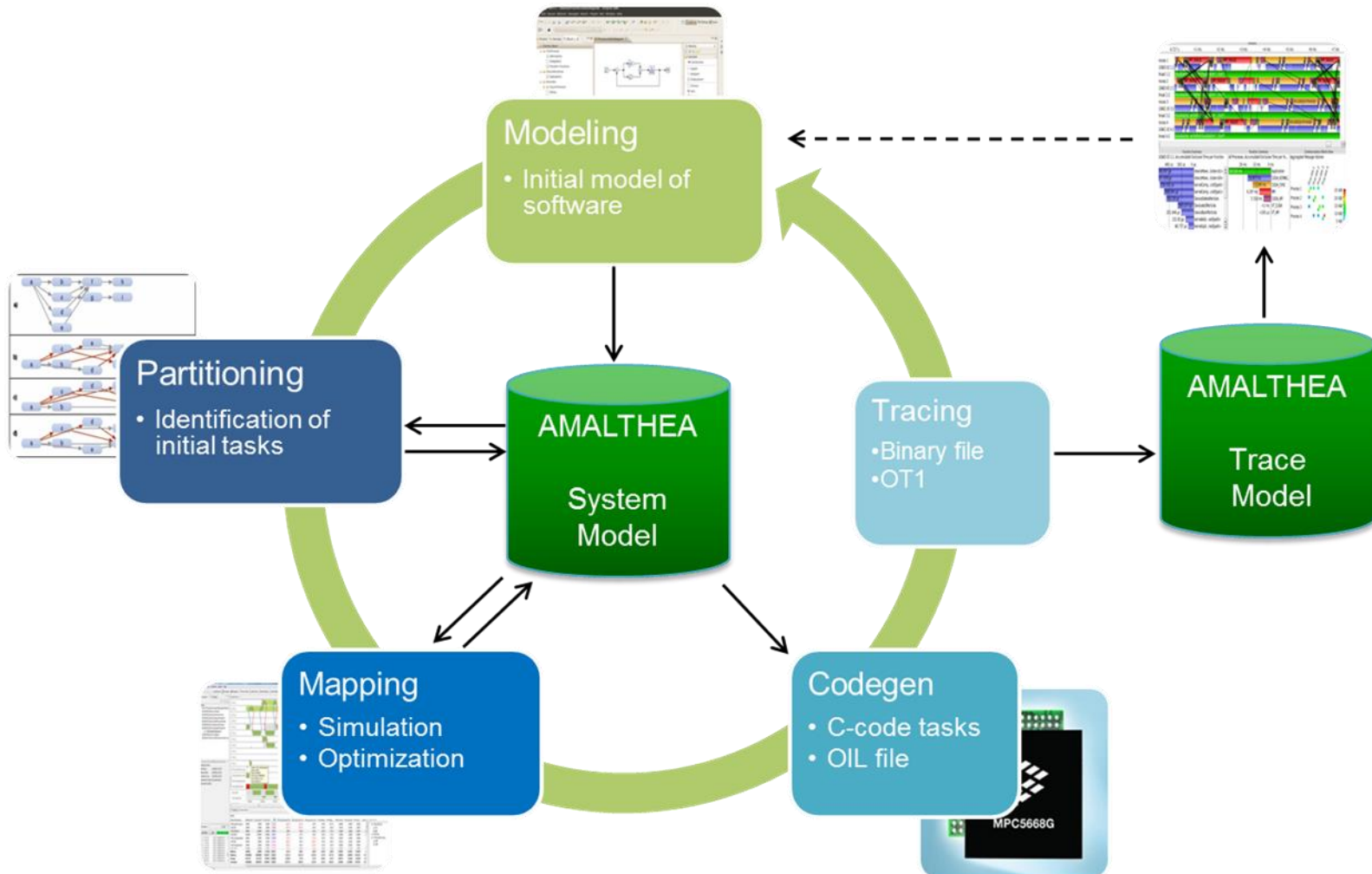
**SysML parametrics**: formulas expressing physical laws (or CPS ones)



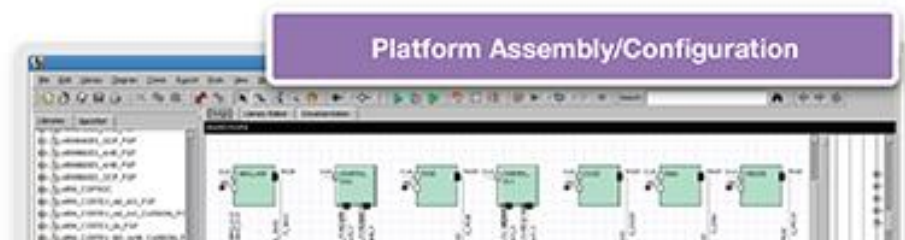
User provides allocations

Analysis tools compute end-to-end latency and other performance measurements

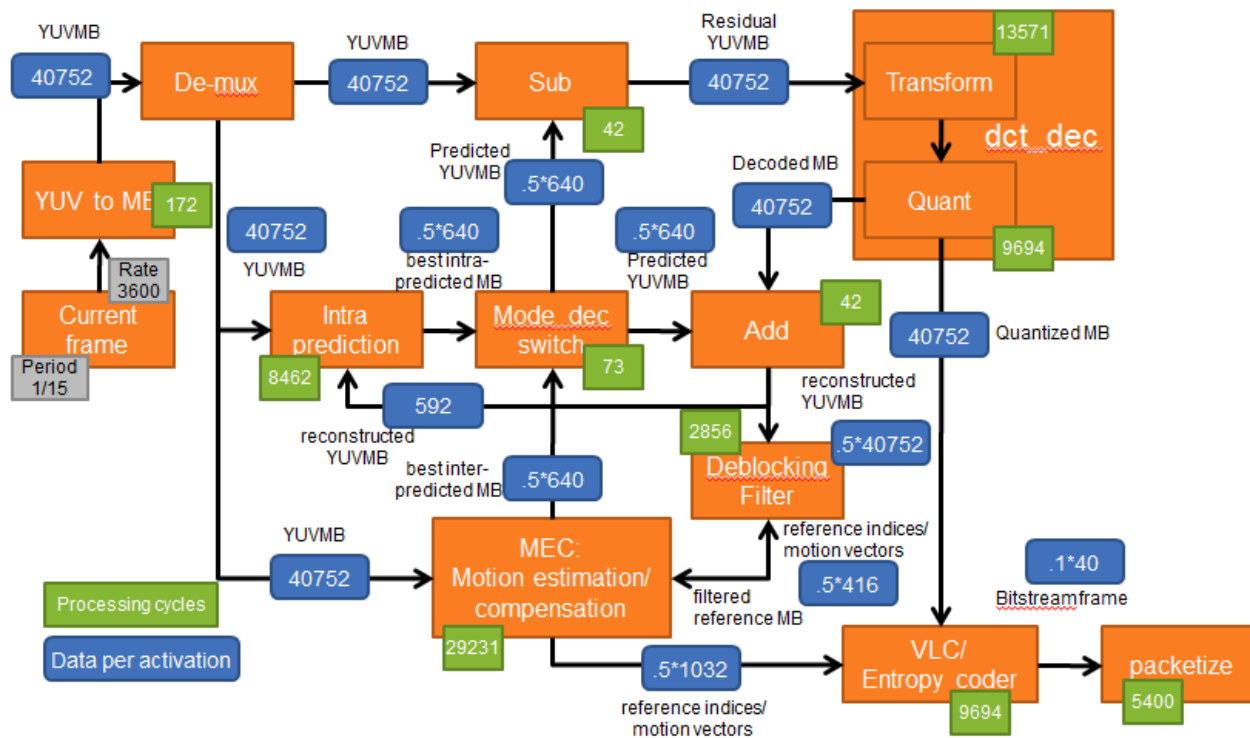
# Example: Amalthea (Bosch et al)







## Example: Synopsys Platform Architect



Many other attempts:

Metropolis  
StreamIt/Raw/Tilera  
SDF3/Aelite  
ForsyDE/Nostrum  
SigmaC/Kalray

...

Function	Self (ns)	Time in Self (ns)	Time in Self/Children (ns)	Callers	Annotations
COMB_VLSI_BasicBlock	10.7	10.7	10.7	1	10
COMB_VLSI_BasicBlock	10.7	10.7	10.7	1	10
COMB_VLSI_BasicBlock	10.7	10.7	10.7	1	10
COMB_VLSI_BasicBlock	10.7	10.7	10.7	1	10
COMB_VLSI_BasicBlock	10.7	10.7	10.7	1	10

Applications and Architectures should be independently described, then (only) fitted together :

- an application may be mapped to multiple execution platforms
- a SW/HW platform is versatile and supports many application

At high-level, simple cost function formulas may be in order, at lower-level a more dynamic simulation relation needed, and if applications have static control the difference is less

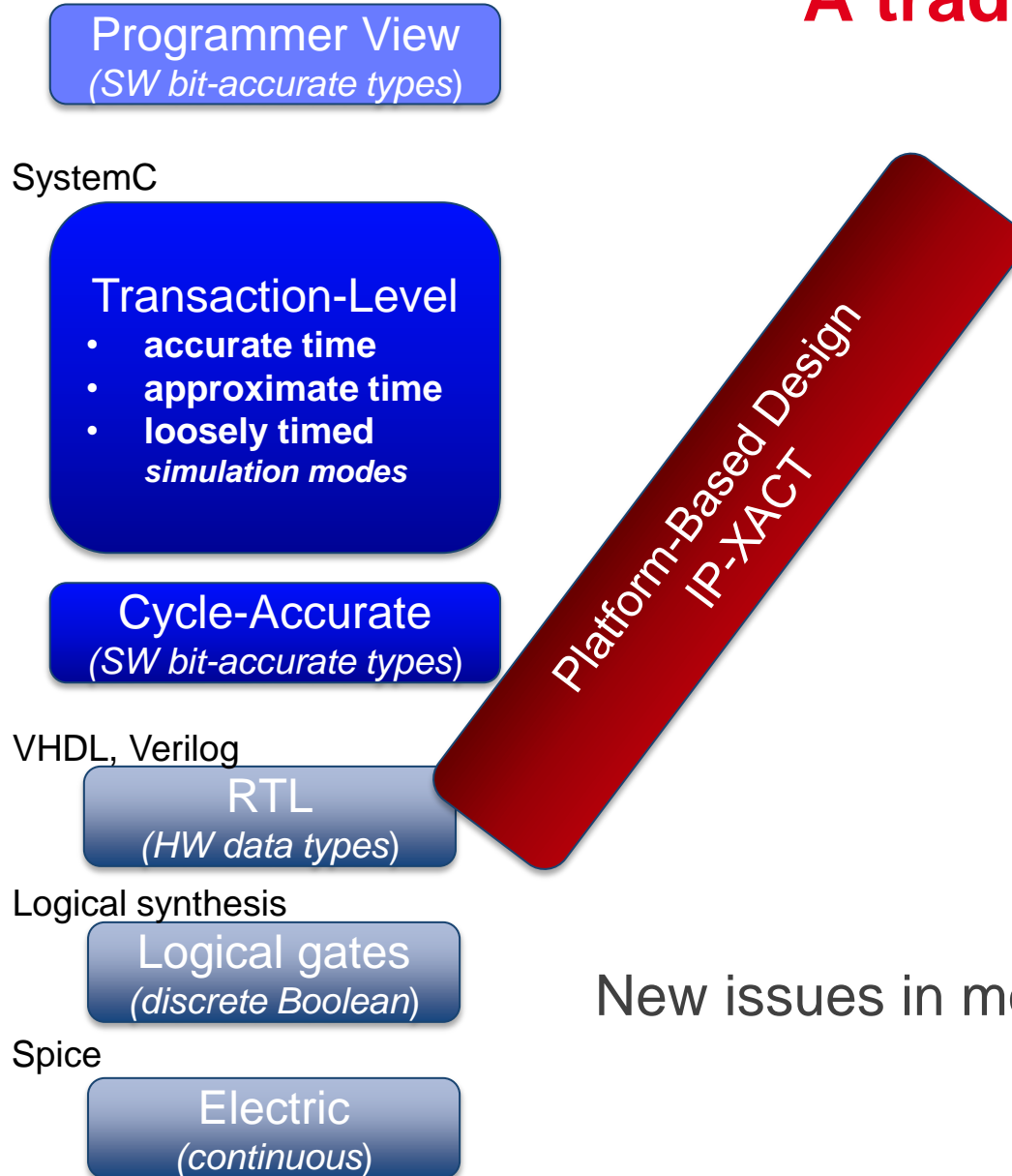
Full control over the (closed) system is assumed. In our case:

**High-Performance or Real-Time Embedded Programming (HPeC)**

# 2

## **Around Hardware Computer-Aided Design** **Sous-titre facultatif**

# A tradition of models

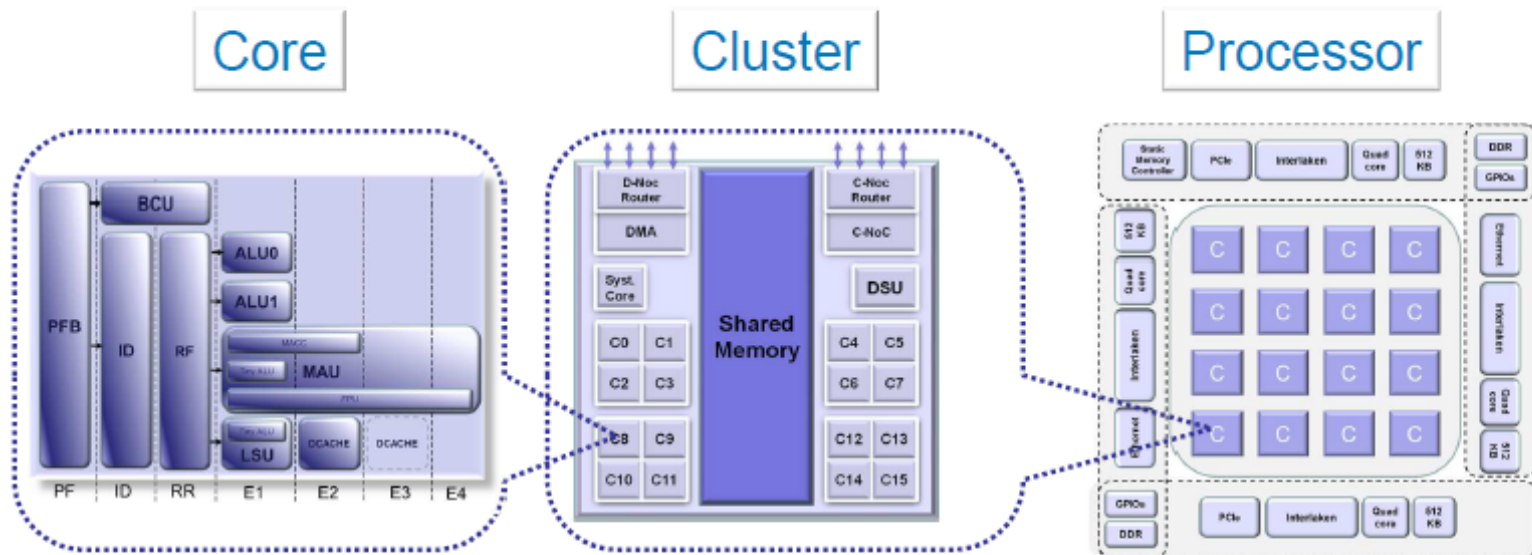


New issues in modeling and simulation

# An example many-core time-deterministic processor



## MPPA<sup>®</sup>-256 Integrated Manycore Processor



- 5-issue VLIW architecture
- Predictability & energy efficiency
- FPU: 32bits / 64 bits IEEE 754
- MMU for rich OS support

- 16 + 1 cores
- NoC Tx and Rx interfaces
- 2 MB of shared memory

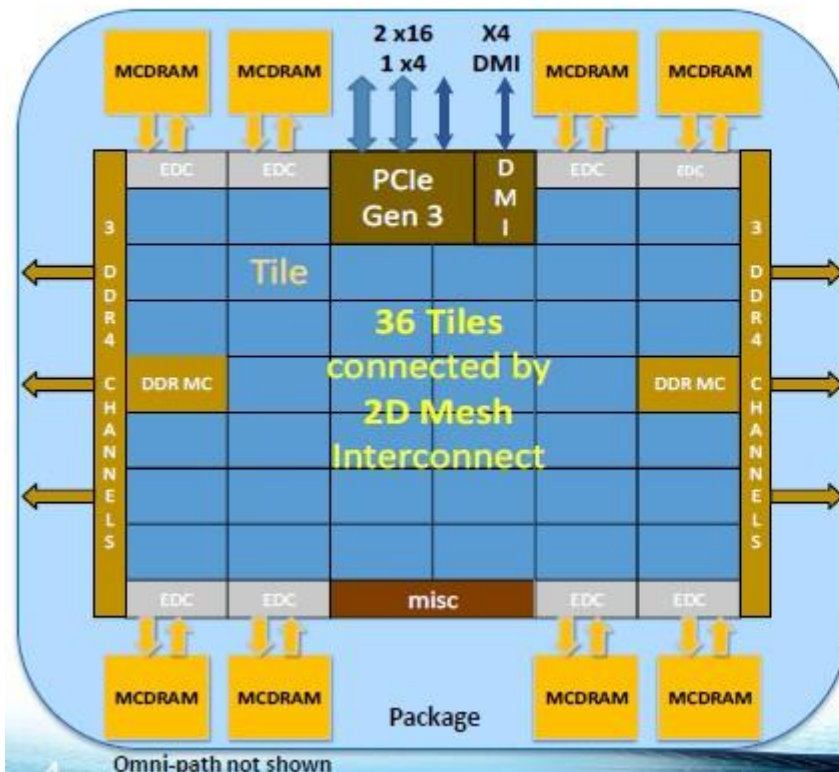
- 2 NoC to connect clusters and I/O
- 47MB of on-chip Memory
- 4x Quad core SMP on each of the four I/O subsystems

### ▪ Execution time predictability

- “Timing Compositional Core”, no timing anomalies (Wilhelm)
- Cluster configuration where each memory bank accessed by 1 core
- Network calculus applied to the Data NoC to bound transfer latency

# A large audience example

## Knights Landing Overview



### TILE

2 VPU	CHA	2 VPU
Core	1MB L2	Core

**Chip: 36 Tiles** interconnected by **2D Mesh**

**Tile: 2 Cores + 2 VPU/core + 1 MB L2**

**Memory: MCDRAM: 16 GB on-package; High BW**

**DDR4: 6 channels @ 2400 up to 384GB**

**IO: 36 lanes PCIe Gen3. 4 lanes of DMI for chipset**

**Node: 1-Socket only**

**Fabric: Omni-Path on-package (not shown)**

**Vector Peak Perf: 3+TF DP and 6+TF SP Flops**

**Scalar Perf: ~3x over Knights Corner**

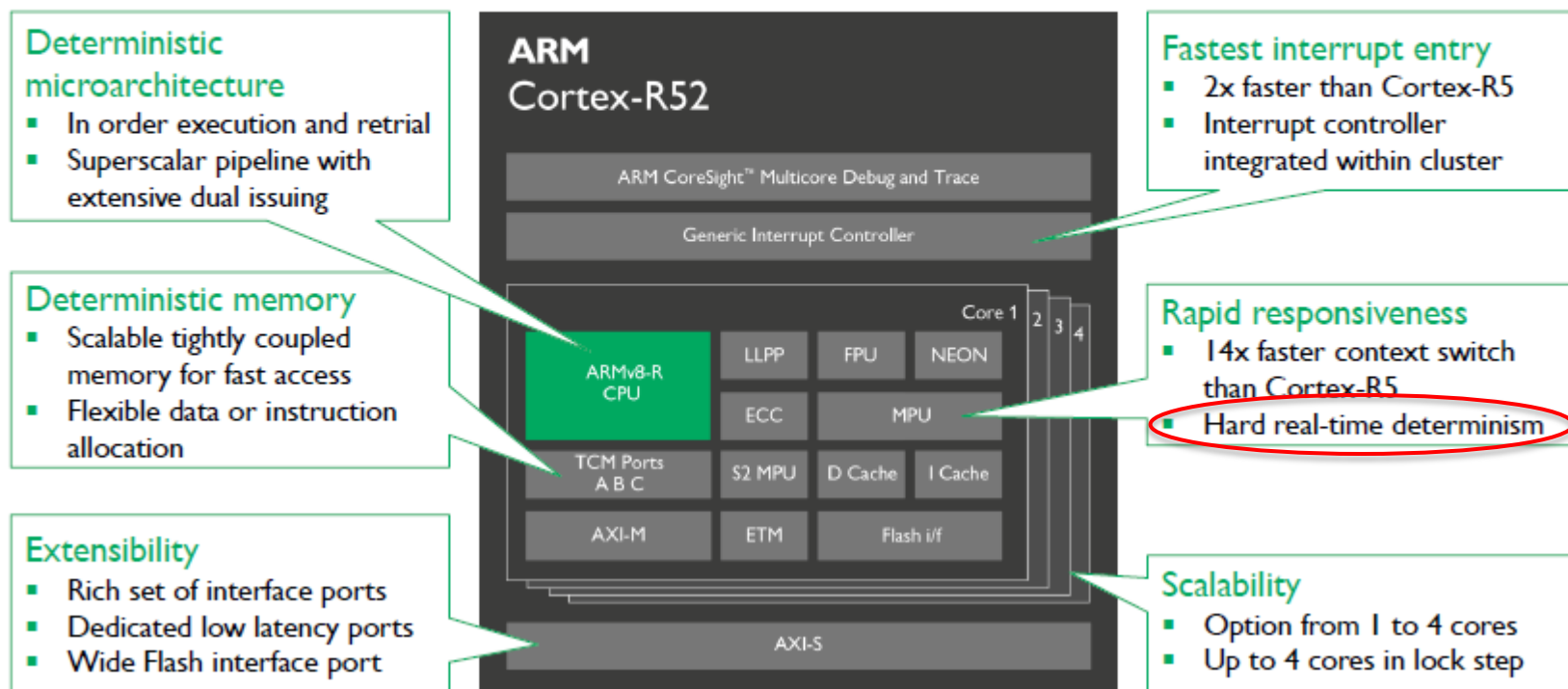
**Streams Triad (GB/s): MCDRAM : 400+; DDR: 90+**

Source Intel: All products, computer systems, dates and figures specified are preliminary based on current expectations, and are subject to change without notice. KNL data are preliminary based on current expectations and are subject to change without notice. 1 Binary Compatible with Intel Xeon processors using Haswell Architecture Set (Intel Xeon Phi). \*Bandwidth numbers are based on STREAM-like memory access pattern when MCDRAM and DDR4 memory. Results have been estimated based on internal Intel analysis and are not intended for commercial purposes only. Any difference in system results or software optimization may affect memory access performance.



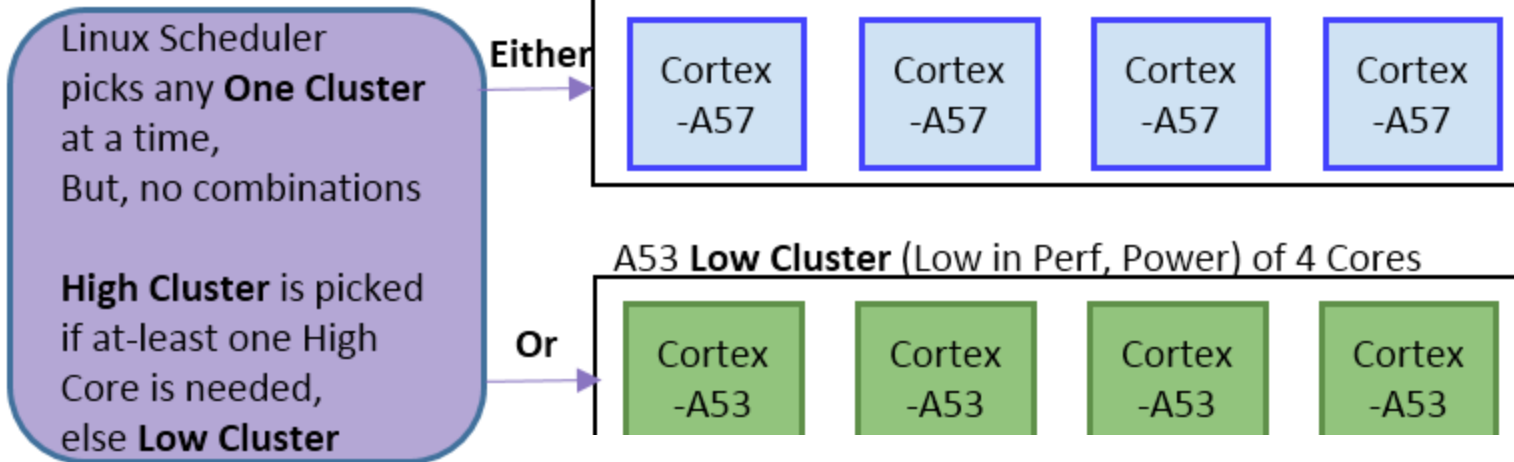
# An example Processor for car applications

## Built for real-time determinism



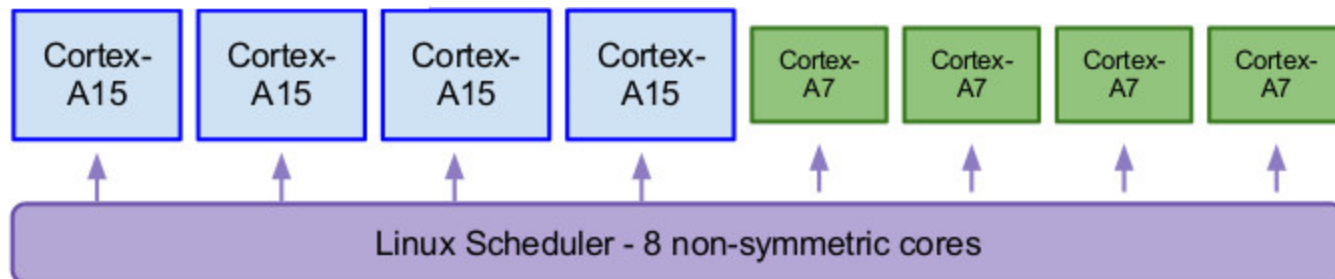
# big. LITTLE alternative mappings

## Clustered Switching:



$$t.duration = \sum_{r \in R} \sum_{o \in O} \underbrace{((t.op == o) \times (t.map == r) \times cost[t][r][o])}_{\text{Boolean decision variables}}$$

## Global Task Scheduling:



**What form of info does one need to organize scheduling ?**  
(cf. upcoming talk by Emilien Kofman)



**Big issues in practice are**

- **Overlapping communications with computations**
- **Dealing with different (physical) clock speeds**  
Processor cores run faster than bus/NoC, than external memory
- **Dealing with performance and low-power (and temperature)**  
Different clock domains and power domains dark silicon, DVFS...

**Ambition is to make the architectural platform model time-predictable**

(as much as possible)

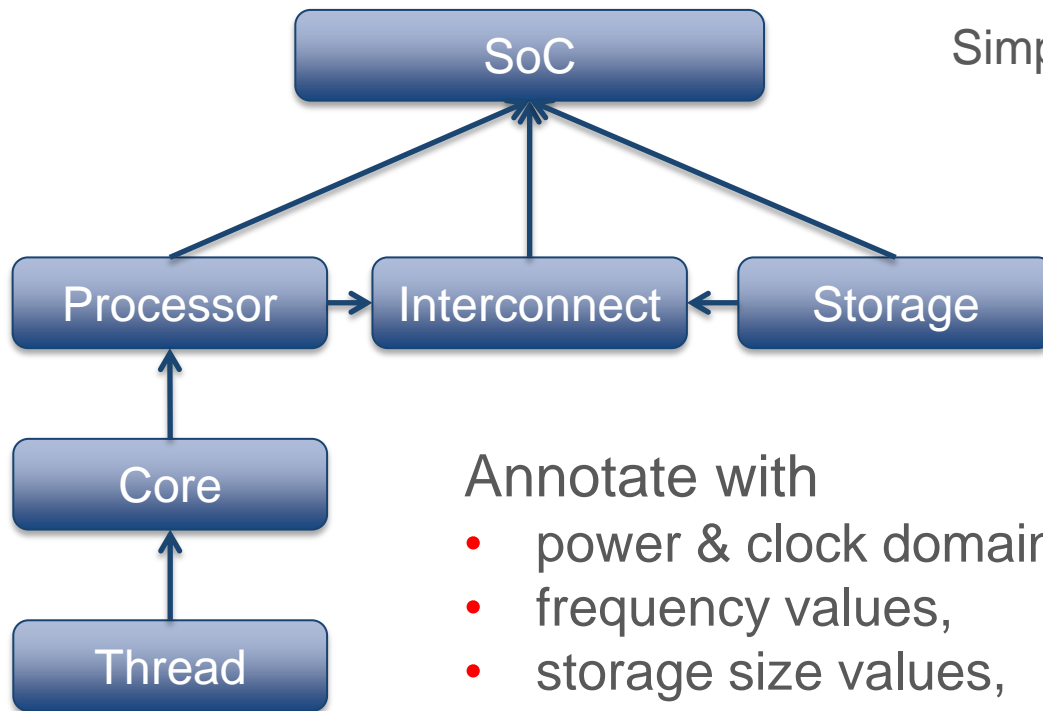
- Either by supposing the real architecture is
- Or by using it in a way that computations and communications fit inside their allotted processing and interconnect resources
  - Avoiding cache misses
  - Avoiding media contention

Requires again static, data-independent control in applications: static planning

**Simple relations between dimensions, here again ?**

(Hardware Architects use Excel for Timing-Closure )...

## Sketch: just enough SoC structure (to support annotation and build constraints)



Simple metamodel ☺ (vs MARTE ☹)

Annotate with

- power & clock domains,
- frequency values,
- storage size values,
- WCET,
- WCCT,
- capacity to overlap computation with communication (DMA and NIC)

# 3

## Around Software Engineering

Programming Models *for parallel processing*

# Parallel Programming Models

## OpenMP (shared memory)

1. Annotations to instruct the compiler on potential parallelism(s):
2. Parallel for-loops, regions, synchro barriers (ingredients of static control programs)
3. With successive versions, always more annotation types (simd, target device)
4. Also more and more ways to inquire the platform about its dimensions and set up affinities  
→ ways to instruct “some” mapping “by hand”

## OpenCL (data parallelism)

1. Programs split as sets of kernels/tasks, to be applied in a data-parallel fashion
2. Same remark as for OpenMP (4) above

## MPI (distributed memory task parallelism and streaming)

1. Networks of (parallel) Processes with message-passing
2. Description of interconnect topology for platform (regular graphs)
3. Placement done by compiler (obscure, may use affinity)

In all cases ,generality (no global static control assumption) is desired.

Some mapping (affinity) and scheduling always done, but without any search for optimality.

No orthogonality (architecture description hinted inside application)

# OpenMP Multiprocessing

- Originally, annotation pragmas to indicate which for-loops could be executed in parallel (if distinct iterations data-independent)
- Supposedl equivalent to sequential form, and shared memory (unlike MPI)
- Issue: no real check that pragmas are correct (remember polyhedral model a few foils ago).
- New in v4: Tasks (so that scheduling is really dynamic); goes again the idea of planning before-hand (compilation)
- More and more pragmas (simd, target...) are letting the actual (supposed) architecture crawl into the program description (no orthogonality) Appli should first be archi agnostic, then only made archi aware by compilation

Execution environment routines affect and monitor threads, processors, and the parallel environment. The library routines are external functions with "C" linkage.

## Execution Environment Routines

### **omp\_set\_num\_threads** [3.2.1] [3.2.1]

Affects the number of threads used for subsequent parallel regions not specifying a `num_threads` clause, by setting the value of the first element of the `nthreads-var` ICV of the current task to `num_threads`.

**void** `omp_set_num_threads(int num_threads);`

### **omp\_get\_num\_threads** [3.2.2] [3.2.2]

Returns the number of threads in the current team. The binding region for an `omp_get_num_threads` region is the innermost enclosing parallel region.

**int** `omp_get_num_threads(void);`

### **omp\_get\_max\_threads** [3.2.3] [3.2.3]

Returns an upper bound on the number of threads that could be used to form a new team if a parallel construct without a `num_threads` clause were encountered after execution returns from this routine.

**int** `omp_get_max_threads(void);`

### **omp\_get\_thread\_num** [3.2.4] [3.2.4]

Returns the thread number of the calling thread within the current team.

**int** `omp_get_thread_num(void);`

### **omp\_get\_num\_procs** [3.2.5] [3.2.5]

Returns the number of processors that are available to the device at the time the routine is called.

**int** `omp_get_num_procs(void);`

### **omp\_in\_parallel** [3.2.6] [3.2.6]

Returns true if the `active-levels-var` ICV is greater than zero; otherwise it returns false.

**int** `omp_in_parallel(void);`

### **omp\_set\_dynamic** [3.2.7] [3.2.7]

Enables or disables dynamic adjustment of the number of threads available for the execution of subsequent parallel regions by setting the value of the `dyn-var` ICV.

**void** `omp_set_dynamic(int dynamic_threads);`

### **omp\_get\_dynamic** [3.2.8] [3.2.8]

This routine returns the value of the `dyn-var` ICV, which is true if dynamic adjustment of the number of threads is enabled for the current task.

**int** `omp_get_dynamic(void);`

### **omp\_get\_cancellation** [3.2.9] [3.2.9]

Returns the value of the `cancel-var` ICV, which is true if cancellation is activated; otherwise it returns false.

**int** `omp_get_cancellation(void);`

### **omp\_set\_nested** [3.2.10] [3.2.10]

Enables or disables nested parallelism, by setting the `nest-var` ICV.

**void** `omp_set_nested(int nested);`

### **omp\_get\_nested** [3.2.11] [3.2.10]

Returns the value of the `nest-var` ICV, which indicates if nested parallelism is enabled or disabled.

**int** `omp_get_nested(void);`

### **omp\_set\_schedule** [3.2.12] [3.2.12]

Affects the schedule that is applied when runtime is used as schedule kind, by setting the value of the `run-sched-var` ICV.

**void** `omp_set_schedule(omp_sched_t kind, int chunk_size);`

**kind:** One of the following, or an implementation-defined schedule:

```
omp_sched_static    = 1
omp_sched_dynamic   = 2
omp_sched_guided    = 3
omp_sched_auto      = 4
```

### **omp\_get\_schedule** [3.2.13] [3.2.13]

Returns the value of `run-sched-var` ICV, which is the schedule applied when runtime schedule is used.

**void** `omp_get_schedule(omp_sched_t *kind, int *chunk_size);`

See **kind** for `omp_set_schedule`.

### **omp\_get\_thread\_limit** [3.2.14] [3.2.14]

Returns the value of the `thread-limit-var` ICV, which is the maximum number of OpenMP threads available.

**int** `omp_get_thread_limit(void);`

### **omp\_set\_max\_active\_levels** [3.2.15] [3.2.15]

Limits the number of nested active parallel regions, by setting `max-active-levels-var` ICV.

**void** `omp_set_max_active_levels(int max_levels);`

### **omp\_get\_max\_active\_levels** [3.2.16] [3.2.16]

Returns the value of `max-active-levels-var` ICV, which determines the maximum number of nested active parallel regions.

**int** `omp_get_max_active_levels(void);`

### **omp\_get\_level** [3.2.17] [3.2.17]

For the enclosing device region, returns the `levels-var` ICV, which is the number of nested parallel regions that enclose the task containing the call.

**int** `omp_get_level(void);`

### **omp\_get\_ancestor\_thread\_num** [3.2.18] [3.2.18]

Returns, for a given nested level of the current thread, the thread number of the ancestor of the current thread.

**int** `omp_get_ancestor_thread_num(int level);`

### **omp\_get\_team\_size** [3.2.19] [3.2.19]

Returns, for a given nested level of the current thread, the size of the thread team to which the ancestor or the current thread belongs.

**int** `omp_get_team_size(int level);`

### **omp\_get\_active\_level** [3.2.20] [3.2.20]

Returns the value of the `active-level-var` ICV, which determines the number of active, nested parallel regions enclosing the task that contains the call.

**int** `omp_get_active_level(void);`

### **omp\_in\_final** [3.2.21] [3.2.21]

Returns true if the routine is executed in a final task region; otherwise, it returns false.

**int** `omp_in_final(void);`

### **omp\_get\_proc\_bind** [3.2.22] [3.2.22]

Returns the thread affinity policy to be used for the subsequent nested parallel regions that do not specify a `proc_bind` clause.

**omp\_proc\_bind\_t** `omp_get_proc_bind(void);`

Returns one of:

```
omp_proc_bind_false = 0
omp_proc_bind_true  = 1
omp_proc_bind_master = 2
omp_proc_bind_close  = 3
omp_proc_bind_spread = 4
```

### **omp\_get\_num\_places** [3.2.23]

Returns the number of places available to the execution environment in the place list.

**int** `omp_get_num_places(void);`

### **omp\_get\_place\_num\_procs** [3.2.24]

Returns the number of processors available to the execution environment in the specified place.

**int** `omp_get_place_num_procs(int place_num);`

### **omp\_get\_place\_proc\_ids** [3.2.25]

Returns the numerical identifiers of the processors available to the execution environment in the specified place.

**void** `omp_get_place_proc_ids(int place_num, int *ids);`

### **omp\_get\_place\_num** [3.2.26]

Returns the place number of the place to which the encountering thread is bound.

**int** `omp_get_place_num(void);`

### **omp\_get\_partition\_num\_places** [3.2.27]

Returns the number of places in the place partition of the innermost implicit task.

**int** `omp_get_partition_num_places(void);`

### **omp\_get\_partition\_place\_nums** [3.2.28]

Returns the list of place numbers corresponding to the places in the `place-partition-var` ICV of the innermost implicit task.

**void** `omp_get_partition_place_nums(int *place_nums);`



**The OpenCL Platform Layer**

The OpenCL platform layer implements platform-specific features that allow applications to query OpenCL devices, device configuration information, and to create OpenCL contexts using one or more devices. Items in blue apply when the appropriate extension is supported.

**Querying Platform Info & Devices [4.1-2] [9.16.9]**

```
cl_int clGetPlatformIDs (cl_uint num_entries,
                        cl_platform_id *platforms, cl_uint *num_platforms)
```

```
cl_int clGetPlatformIDsKHR (cl_uint num_entries,
                           cl_platform_id *platforms, cl_uint *num_platforms)
```

```
cl_int clGetPlatformInfo (cl_platform_id platform,
                        cl_platform_info param_name,
                        size_t param_value_size, void *param_value,
                        size_t *param_value_size_ret)
```

param\_name: CL\_PLATFORM\_PROFILE, VERSION},  
CL\_PLATFORM\_NAME, VENDOR, EXTENSIONS},  
CL\_PLATFORM\_ICD\_SUFFIX\_KHR [Table 4.1]

```
cl_int clGetDeviceIDs (cl_platform_id platform,
                      cl_device_type device_type, cl_uint num_entries,
                      cl_device_id *devices, cl_uint *num_devices)
```

device\_type: [Table 4.2]  
CL\_DEVICE\_TYPE\_ACCELERATOR, ALL, CPU},  
CL\_DEVICE\_TYPE\_CUSTOM, DEFAULT, GPU}

```
cl_int clGetDeviceInfo (cl_device_id device,
                      cl_device_info param_name,
                      size_t param_value_size, void *param_value,
                      size_t *param_value_size_ret)
```

param\_name: [Table 4.3]  
CL\_DEVICE\_ADDRESS\_BITS, CL\_DEVICE\_AVAILABLE,  
CL\_DEVICE\_BUILT\_IN\_KERNELS,  
CL\_DEVICE\_COMPILER\_AVAILABLE,  
CL\_DEVICE\_DOUBLE, HALF, SINGLE}, FP\_CONFIG,  
CL\_DEVICE\_ENDIAN\_LITTLE, CL\_DEVICE\_EXTENSIONS,  
CL\_DEVICE\_ERROR\_CORRECTION\_SUPPORT,  
CL\_DEVICE\_EXECUTION\_CAPABILITIES,  
CL\_DEVICE\_GLOBAL\_MEM\_CACHE\_SIZE, TYPE},  
CL\_DEVICE\_GLOBAL\_MEM\_CACHELINE\_SIZE, SIZE},  
CL\_DEVICE\_GLOBAL\_VARIABLE\_PREFERRED\_TOTAL\_SIZE,  
CL\_DEVICE\_PREFERRED\_PLATFORM\_LOCAL,  
GLOBAL\_ATOMIC\_ALIGNMENT,  
CL\_DEVICE\_GLOBAL\_VARIABLE\_SHARING,  
CL\_DEVICE\_HOST\_UNIFIED\_MEMORY,  
CL\_DEVICE\_IMAGE\_MAX\_ARRAY\_BUFFER\_SIZE,  
CL\_DEVICE\_IMAGE\_SUPPORT,  
CL\_DEVICE\_IMAGE2D\_MAX\_WIDTH\_HEIGHT},  
CL\_DEVICE\_IMAGE3D\_MAX\_WIDTH\_HEIGHT\_DEPTH},  
CL\_DEVICE\_IMAGE\_BASE\_ADDRESS\_ALIGNMENT,  
CL\_DEVICE\_IMAGE\_PITCH\_ALIGNMENT,  
CL\_DEVICE\_LINKER\_AVAILABLE,  
CL\_DEVICE\_LOCAL\_MEM\_TYPE, SIZE},  
CL\_DEVICE\_MAX\_READ\_IMAGE\_ARGS,  
CL\_DEVICE\_MAX\_WRITE\_IMAGE\_ARGS,  
CL\_DEVICE\_MAX\_CLOCK\_FREQUENCY\_PIPE\_ARGS},  
CL\_DEVICE\_MAX\_COMPUTE\_UNITS\_SAMPLERS},  
CL\_DEVICE\_MAX\_CONSTANT\_ARGS\_BUFFER\_SIZE},  
CL\_DEVICE\_MAX\_MEM\_ALLOC\_PARAMETER\_SIZE,  
CL\_DEVICE\_MAX\_GLOBAL\_VARIABLE\_SIZE,  
CL\_DEVICE\_MAX\_ON\_DEVICE\_QUEUES\_EVENTS},  
CL\_DEVICE\_MAX\_WORK\_GROUP\_SIZE,  
CL\_DEVICE\_MAX\_WORK\_ITEM\_DIMENSIONS\_SIZES},  
CL\_DEVICE\_MEM\_BASE\_ADDR\_ALIGN,  
CL\_DEVICE\_NAME,  
CL\_DEVICE\_NATIVE\_VECTOR\_WIDTH\_CHAR, INT},  
CL\_DEVICE\_NATIVE\_VECTOR\_WIDTH\_LONG, SHORT},  
CL\_DEVICE\_NATIVE\_VECTOR\_WIDTH\_DOUBLE, HALF},  
CL\_DEVICE\_NATIVE\_VECTOR\_WIDTH\_FLOAT,  
CL\_DEVICE\_OPENCL\_C\_VERSION\_PARENT\_DEVICE},  
CL\_DEVICE\_PARTITION\_AFFINITY\_DOMAIN,  
CL\_DEVICE\_PARTITION\_MAX\_SUB\_DEVICES,  
CL\_DEVICE\_PARTITION\_PROPERTIES\_TYPE},  
CL\_DEVICE\_PIPE\_MAX\_ACTIVE\_RESERVATIONS,  
CL\_DEVICE\_PIPE\_MAX\_PACKET\_SIZE,

CL\_DEVICE\_PLATFORM, PRINTF\_BUFFER\_SIZE},  
CL\_DEVICE\_PREFERRED\_VECTOR\_WIDTH\_CHAR, INT},  
CL\_DEVICE\_PREFERRED\_VECTOR\_WIDTH\_DOUBLE,  
CL\_DEVICE\_PREFERRED\_VECTOR\_WIDTH\_HALF,  
CL\_DEVICE\_PREFERRED\_VECTOR\_WIDTH\_LONG,  
CL\_DEVICE\_PREFERRED\_VECTOR\_WIDTH\_SHORT,  
CL\_DEVICE\_PREFERRED\_VECTOR\_WIDTH\_FLOAT,  
CL\_DEVICE\_PREFERRED\_INTEROP\_USER\_SYNC,  
CL\_DEVICE\_PROFILE,  
CL\_DEVICE\_PROFILING\_TIMER\_RESOLUTION,  
CL\_DEVICE\_SPIR\_VERSIONS,  
CL\_DEVICE\_QUEUE\_ON\_DEVICE\_PROPERTIES,  
CL\_DEVICE\_QUEUE\_ON\_HOST\_PROPERTIES,  
CL\_DEVICE\_QUEUE\_ON\_DEVICE\_MAX\_SIZE,  
CL\_DEVICE\_QUEUE\_ON\_DEVICE\_PREFERRED\_SIZE,  
CL\_DEVICE\_REFERENCE\_COUNT\_VENDOR\_ID},  
CL\_DEVICE\_SVM\_CAPABILITIES,  
CL\_DEVICE\_TERMINATE\_CAPABILITY\_KHR,  
CL\_DEVICE\_TYPE\_VENDOR},  
CL\_DEVICE\_DRIVER\_VERSION

**Partitioning a Device [4.3]**

```
cl_int clCreateSubDevices (cl_device_id in_device,
                          const cl_device_partition_property *properties,
                          cl_uint num_devices, cl_device_id *out_devices,
                          cl_uint *num_devices_ret)
```

properties: CL\_DEVICE\_PARTITION\_EQUALLY,  
CL\_DEVICE\_PARTITION\_BY\_COUNTS,  
CL\_DEVICE\_PARTITION\_BY\_AFFINITY\_DOMAIN

```
cl_int clRetainDevice (cl_device_id device)
```

```
cl_int clReleaseDevice (cl_device_id device)
```

**Contexts [4.4]**

```
cl_context clCreateContext (
    const cl_context_properties *properties,
    cl_uint num_devices, const cl_device_id *devices,
    void (CL_CALLBACK *pfn_notify)
    (const char *errinfo, const void *private_info,
     size_t cb, void *user_data),
    void *user_data, cl_int *errcode_ret)
```

**The OpenCL Runtime**

API calls that manage OpenCL objects such as command-queues, memory objects, program objects, kernel objects for kernel functions in a program and calls that allow you to enqueue commands to a command-queue such as executing a kernel, reading, or writing a memory object.

**Command Queues [5.1]**

```
cl_command_queue
clCreateCommandQueueWithProperties (
    cl_context context, cl_device_id device,
    const cl_command_queue_properties *properties,
    cl_int *errcode_ret)
```

properties: [Table 5.1] CL\_QUEUE\_SIZE,  
CL\_QUEUE\_PROPERTIES (bitfield which may be set to an OR of CL\_QUEUE\_\* where \* may be: OUT\_OF\_ORDER\_EXEC\_MODE\_ENABLE, PROFILING\_ENABLE, ON\_DEVICE\_DEFAULT)

```
cl_int clRetainCommandQueue (
    cl_command_queue command_queue)
```

```
cl_int clReleaseCommandQueue (
    cl_command_queue command_queue)
```

```
cl_int clGetCommandQueueInfo (
    cl_command_queue command_queue,
    cl_command_queue_info param_name,
    size_t param_value_size, void *param_value,
    size_t *param_value_size_ret)
```

param\_name: [Table 5.2] CL\_QUEUE\_CONTEXT,  
CL\_QUEUE\_DEVICE, CL\_QUEUE\_SIZE,  
CL\_QUEUE\_REFERENCE\_COUNT,  
CL\_QUEUE\_PROPERTIES

# OpenCL

and explicit data transfer for GPU

verified the processing elements, and  
have arrived (finger-crossed)  
computation and communication

# Message-Passing Interface (MPI)

- Highly used for High-Performance Computing (grids, clusters)
- Mostly message-passing and synchronization primitives, to be applied with any general-purpose sequential language (in fact Fortran or C++ generally)
- Point-to-point or collective (broadcast, all-to-all) communications
- Means to describe virtual network topology as graph
- Big issue is to assign processes to processors (statically?), and to map the virtual communication onto the real interconnect infrastructures
- Both prominent frameworks (OpenMPI and MPICH) have fancy dedicated library, commercial offers are made by specific hardware vendors
- Mapping done at runtime, but what about static control processes and static mapping at compile time ?



## •Communicators with Topology

## MPI quick card (excerpt)

- Create with cartesian topology. ( § 6.5.1)
- int **MPI\_Cart\_create** (MPI\_Comm comm\_old,  
•int ndims, int \*dims, int \*periods, int  
•reorder, MPI\_Comm \*comm\_cart)
- Suggest balanced dimension ranges. ( § 6.5.2)
- int **MPI\_Dims\_create** (int nnodes, int  
•ndims, int \*dims)
- Determine rank from cartesian coordinates. ( § 6.5.4)
- int **MPI\_Cart\_rank** (MPI\_Comm comm, int  
•\*coords, int \*rank)
- Determine cartesian coordinates from rank. ( § 6.5.4)
- int **MPI\_Cart\_coords** (MPI\_Comm comm, int  
•rank, int maxdims, int \*coords)
- Determine ranks for cartesian shift. ( § 6.5.5)
- int **MPI\_Cart\_shift** (MPI\_Comm comm, int  
•direction, int disp, int \*rank\_source,  
•int \*rank\_dest)
- Split into lower dimensional sub-grids. ( § 6.5.6)
- int **MPI\_Cart\_sub** (MPI\_Comm comm, int  
•\*remain\_dims, MPI\_Comm \*newcomm)
- Related Functions:* MPI\_Graph\_create, MPI\_Topo\_test,  
•MPI\_Graphdims\_get, MPI\_Graph\_get,  
•MPI\_Cartdim\_get, MPI\_Cart\_get,  
•MPI\_Graph\_neighbors\_count, MPI\_Graph\_neighbors,  
•MPI\_Cart\_map, MPI\_Graph\_map

# Compilation and Runtime Execution

- HPC Runtime: StarPU, XKaapi, RMC
- Parallel compilation & Polyhedral model : ClooG, Graphite, R-Stream
  - (often produces OpenMP-ish code)

These always operate by trying to exhibit maximal parallelism on the application side (and the kind of parallelism wanted first), then try to adjust it dynamically to a “phantom” architecture only reflected by a few features showing availability of simple resources

**Can we do better with a real architecture model?**

(Of course static control restrictions for compile-time decisions)

# Data/Task parallel levels: Nested For-Loop programs with affine bounds

Both code and model:

- Strict separation between indexes and other variables  
(only indexes impact control)
- Other variables are array locations  
(referenced by affine expressions of indexes)

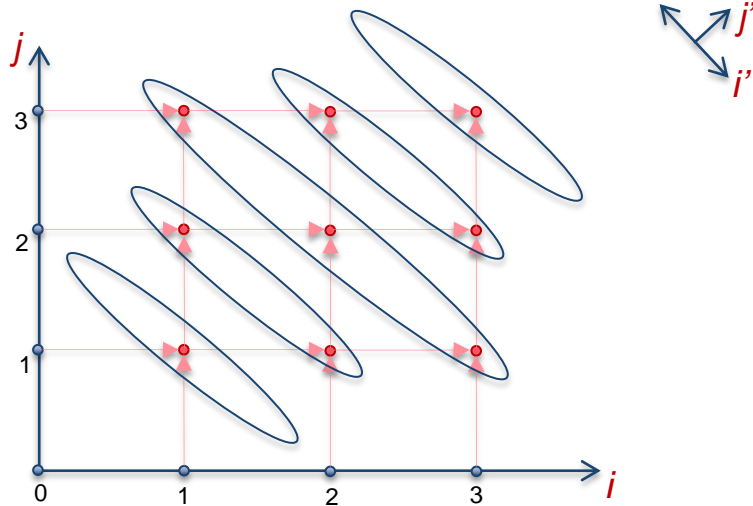
May produce explicit control/data flowgraphs to extracy potential parallelism  
(from dependences between operations)

- Extended: precise but untractable
- Reduced: quality of solution depends on information preserved
  - ☐ Dependence levels (Allen-Kennedy)
  - ☐ Direction vectors (Lamport, Wolf & Lam)
  - ☐ Polyhedral models (Feautrier)

**In all 3 cases, solutions found are expressed by regular (affine) scheduling relations between operations : logical clocks !**

# Silly simple example

```
for j = 1 to N
  for i = 1 to N
    a(i+1,j+1)= a(i, j+1) + a(i+1,j)
  end
end
```

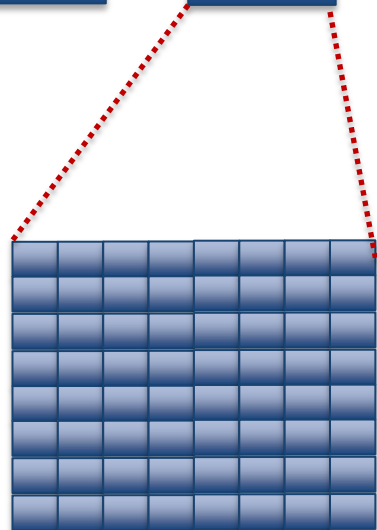


dependences

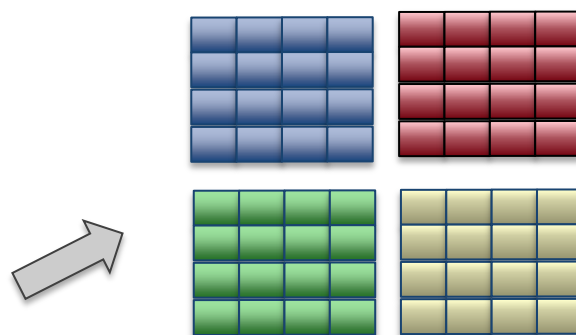
```
for j' = 2 to 2N
  parfor i' = max(1, j'-N) to min(N, j'-1)
    a(i, j'-i')= a(i'-1, j'-i') + a(i', j'-i'-1)
  end
end
```

step(a(i,j))= i+j-1

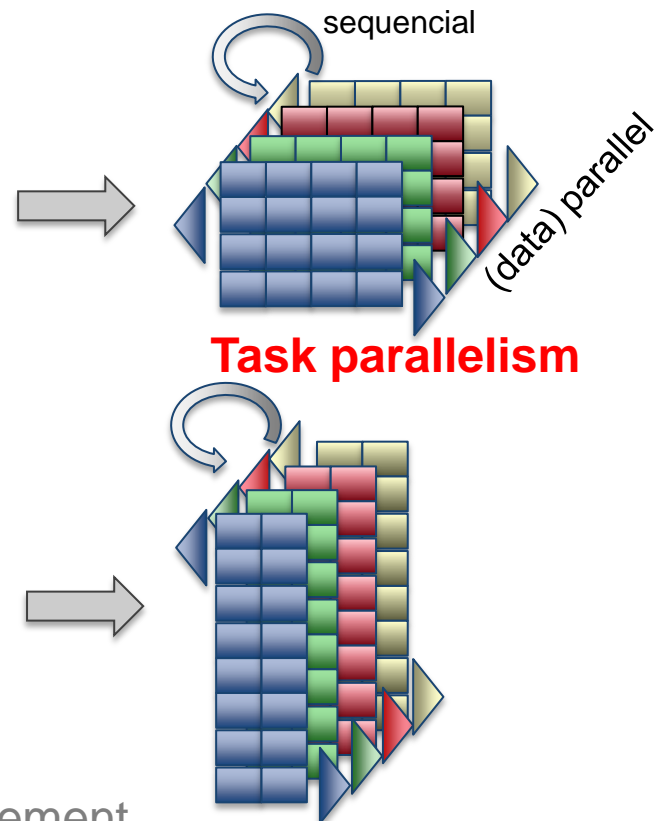
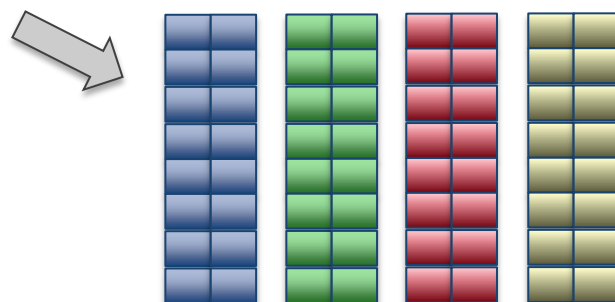
## Task streaming, software pipelining



## Data parallelism



or



## Task parallelism

- Change the granularity
- ...to adjust it to match the size of a single Processing Element
- Locality and shared memory space makes performance highly predictable

# Typical (and useful) library kernels

- **Linear algebra:** LinPack ScalaPack (Blast),...
- **Convolutions** Fast Fourier Transform: FFTW, **Spiral**
- **Neural Networks** (only execution, not training part): a mix of both
- **Signal Processing** (filters) : **Halide**
- Stencil algorithms for numerical analysis ./ scientific computations

This type of effort lead to **Domain-Specific Languages**

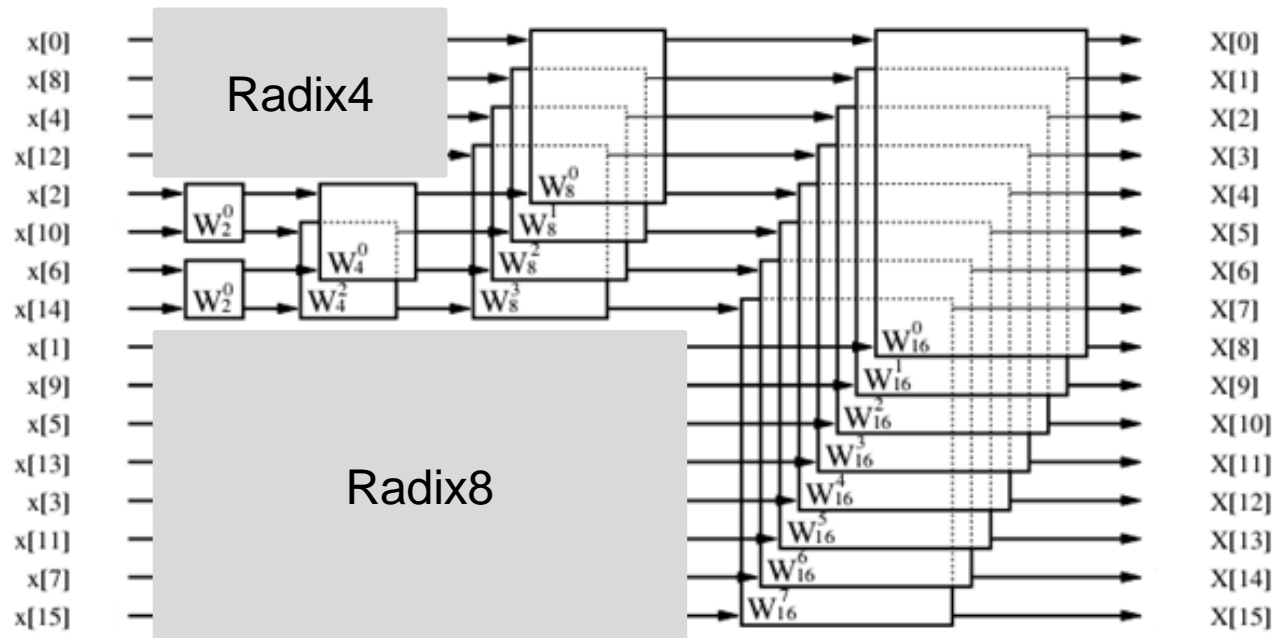
Restricted expressivity: linear or DAG filter streams

Transformations to exhibit some form of parallelism (data or vectorize)

Scheduling often not fully automatic, or obtained by progressive selection (Dynamic Programming, Auto-tuning, or “Auto-Scheduler”...)

...although algorithm stream definitions are often recursive in terms of the data (multi-array) size, so that with a simple core architecture target model it is easy to compute which size fits in (and then loop sequentially to get results). (Application ground floor resembles architecture groundfloor)

# Fast Fourier Transform



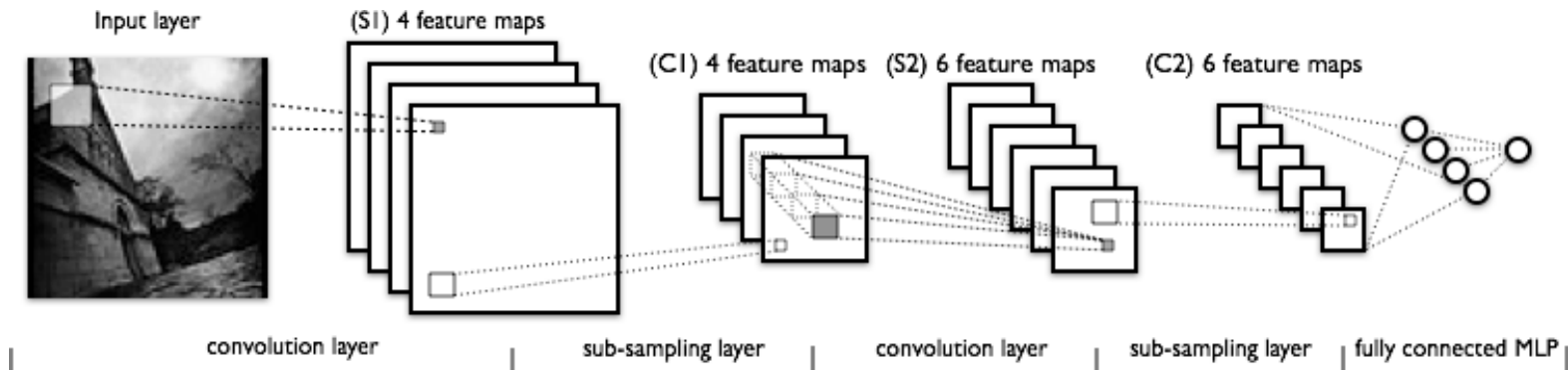
Can easily be decomposed (recursively up to intermediate permutations) along sizes  $2^N$

Given a target processing core (or thread):

- its local register count  $2^R$
- its local memory  $2^M$
- its number of vectorized (simd) or parallel (GPU) ALUs  $2^C$

one can compute by proportions a pattern adjusting the available parallelism (and sequential loops around)

# Convolutional Neural Networks



source: <http://deeplearning.net/tutorial/lenet.html>

Task streamin for layers

+ potential data parallelism according to locality and tile overlap



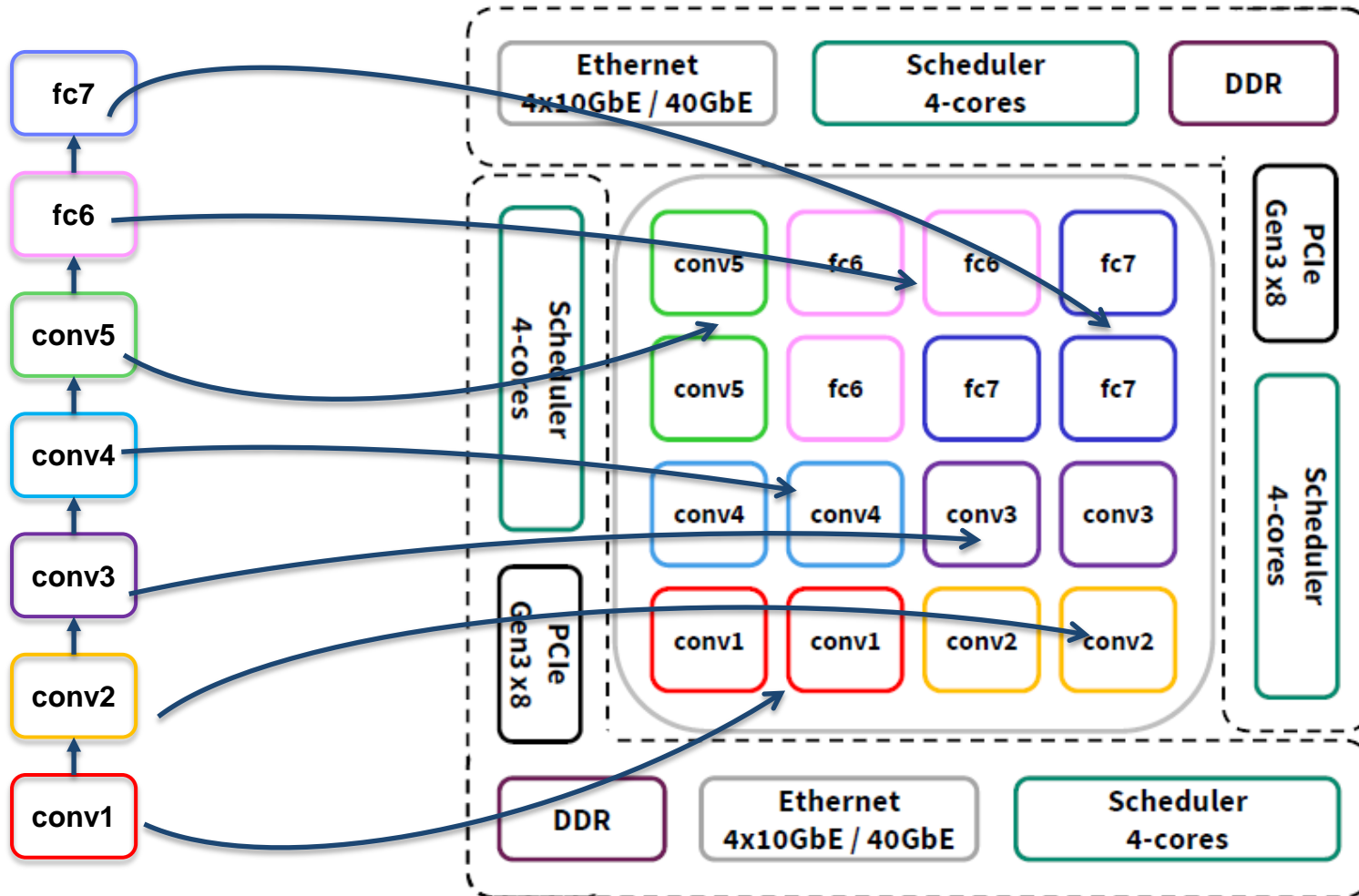
## When

- Application restricted to simple filter DAG/pipeline kernel
- Data value organization well understood (and statically known)
- Simple processor targeted (single core, processing element, or thread)
- Shared memory, low-cost communication (data transfer)
- Mutual correspondance of sizing between appli and archi, and simple cost function formulas,

Then the ground floor mapping of library kernels to core should deliver time-predictable results (WCET)

Then one could consider optimizing full (static control) applications on full platform, designing time-predicatble communication pattterns.

# ConvNets on Kalray MPPA



source: B. Ganne (Kalray) , Machine Learning on Multicores, NeuroStic Days 2015 (in French)

## SIMD CPU cores (4 cores)

- each core enough L1/register memory to execute 4 stages at once :  
*2 256-bit SIMD ALUs, 16 256-bit registers, 128 floating-point numbers*
- FFT 1064 (radix 10, 10 stages) performed as iteration 4+4+2
- in C++ with simd intrinsics

## GPU Processing Elements (40)

- each PE enough L1/register memory to execute 3 stages at one  
*2 128-bit ALUs, 2K bit registers, 64 floating-point numbers*
- FFT 1024 executed as iteration 3+3+3+2
- in OpenCL

To force allocation one needs affinity (dismiss hyperthreading etc), currently not obvious

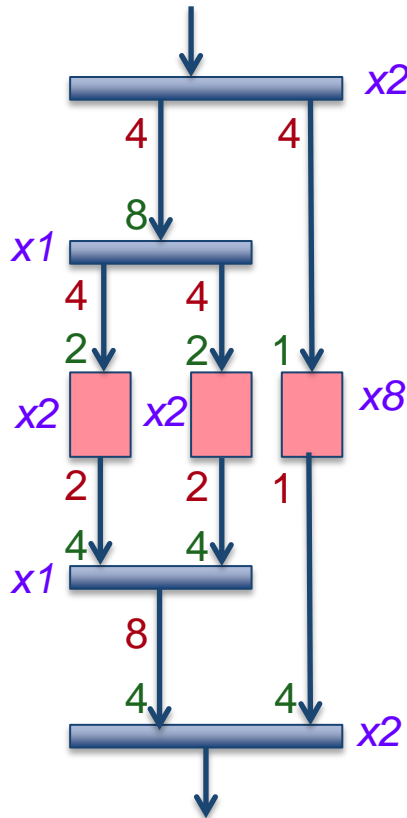
In the end one gets a libray of mapped functions with reasonably accurate cost, to be used by upper floor analyses

# 4

## Concurrent Models of Computation and Communication

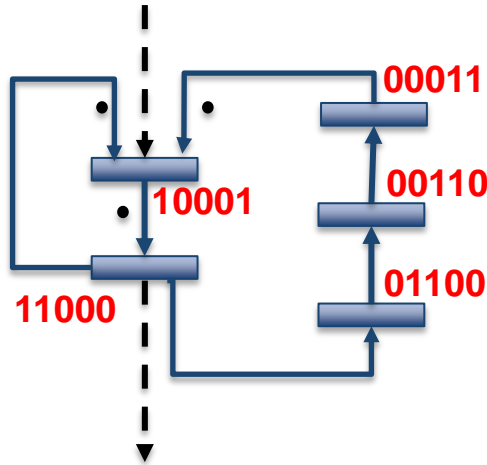
Process Networks, Scheduling and Routing

## SDF (synchronous data-flow process networks)



- Weighted Marked Graphs (conflict-free subset of Petri Nets)
- Distinct issues for acyclic graphs and strongly connected components
- First-level self-timed semantics
- Second stage: optimal scheduling exist, assuming no parallelism limitations
- Static schedules lead to regular form (the binary word of activations of any node is of specific form)
- all nodes adopt the same rate

# Regular schedules in SDF



The noticeable thing here:

- (static) schedules can be represented as regular infinite binary words

$u.(v)$  ,  $u$  initial,  $v$  repeated

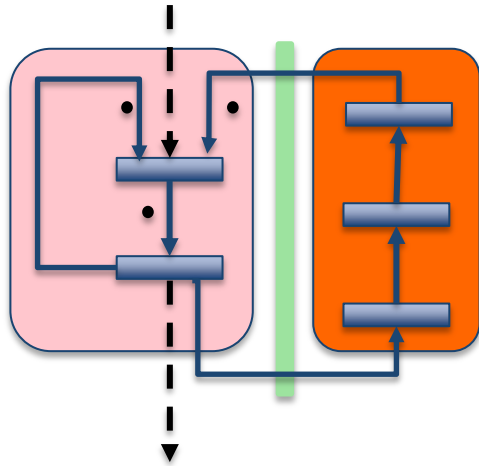
where

- ☐ 1 at location  $n$  means active at step  $n$
- ☐ 0 at location  $n$  means idle at step  $n$

While very primitive, it corresponds to the way ASAP schedule assignment goes  
Later it can be turned into Gantt Charts, when sequences get long

No architectural resource constraints considered here: ideal parallelism

# Regular schedules in SDF

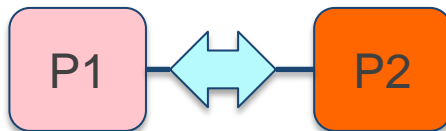


Now what if we map (here superpose) to an execution platform ?

- Cost of communication is no longer uniform (nuMA)
- Computations or communications sharing a single resource may need to be serialized (multitasking)

But the principles of asap scheduling stays the same !

Even more constraints may be applied and task fission/fusion also...



**To what extent can we represent in this abstract setting the phenomena encountered earlier ?**

## Adding predictable deterministic control to SDF

- add a switch node), but with internal switching condition (as in Kahn Process Networks → functional determinism, latency-insensitive design
- done in Cyclo-Static Data-Flow and StreamIt graphs (but without initialization patterns
- Our own view: add regular **switching/routing** patterns in the same flavor as schedule words activating conditions

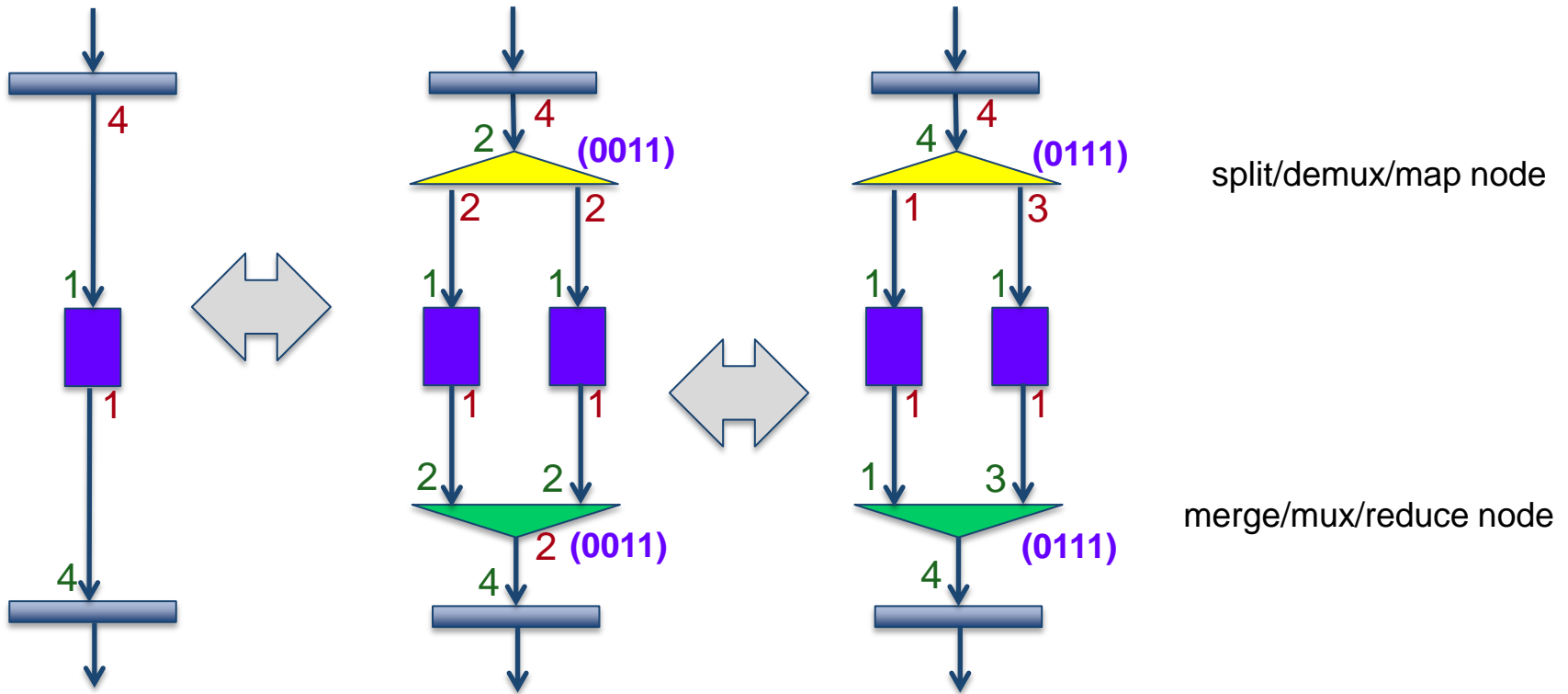
→ KRG process network model

- Study (equational, algebraic) graph transformations that preserve functionality (self-timed)
- these transformations made to change the buffer/data dependencies to adjust to a given platform communication topology graph (eg, NoC)

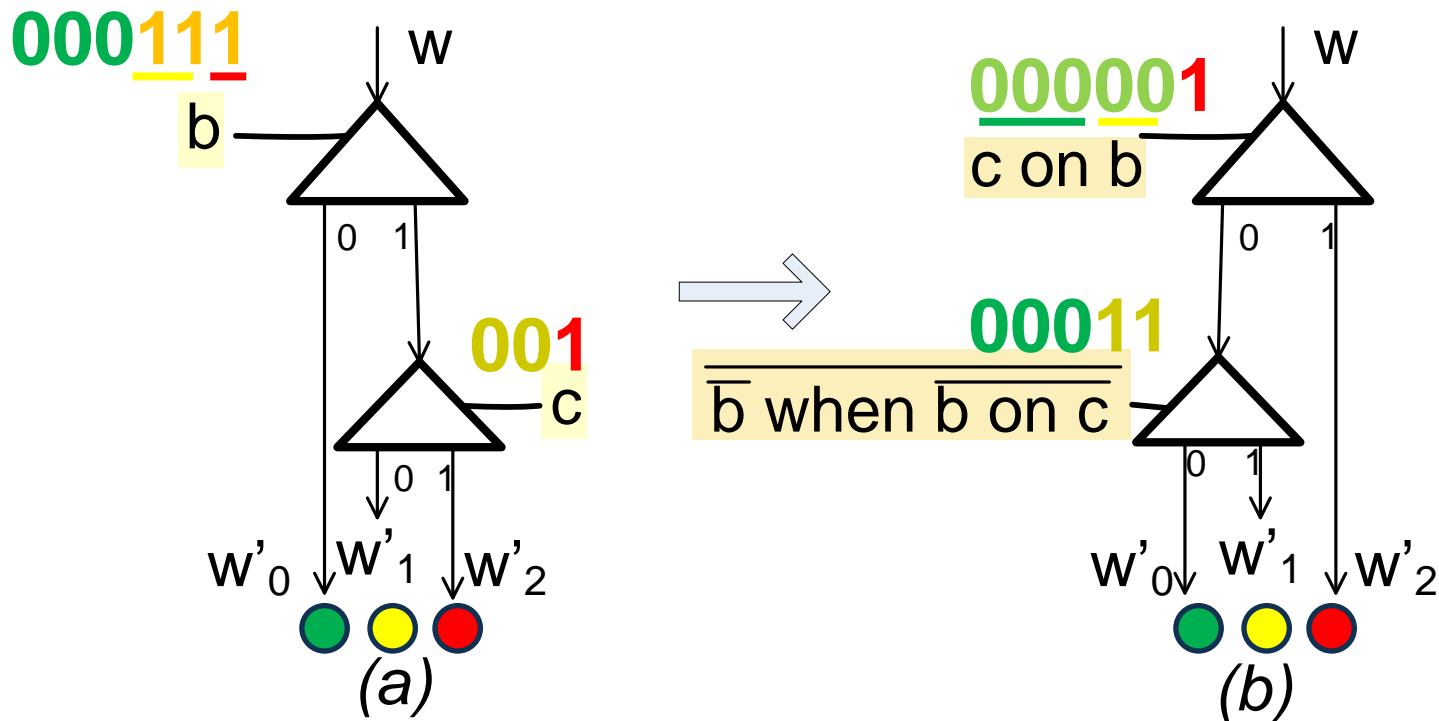




# Silly little KRG example



# Transitivity of Selects



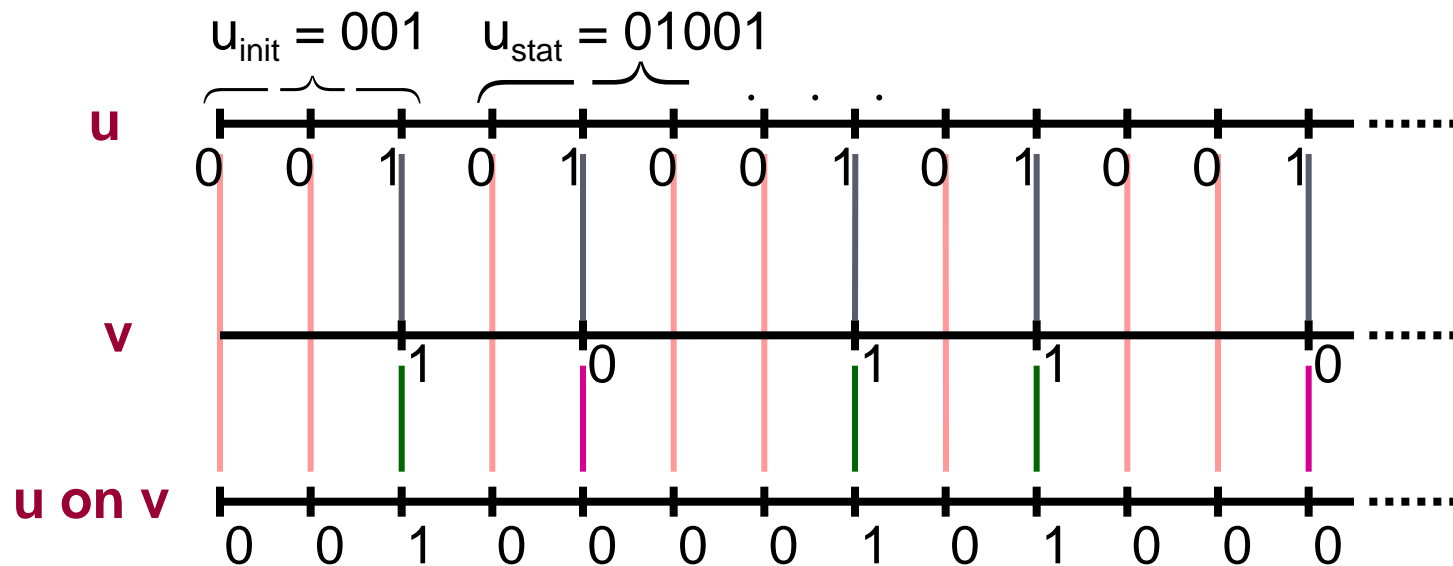
# on/when operators

$$(0.u) \text{ on } v = 0.(u \text{ on } v)$$
$$(1.u) \text{ on } x.v = x.(u \text{ on } v)$$

$$(x.u) \text{ when } (0.v) = u \text{ when } v$$
$$(x.u) \text{ when } (1.v) = x.(u \text{ when } v)$$

# on effects

$$(0.u) \text{ on } v = 0.(u \text{ on } v)$$
$$(1.u) \text{ on } x.v = x.(u \text{ on } v)$$

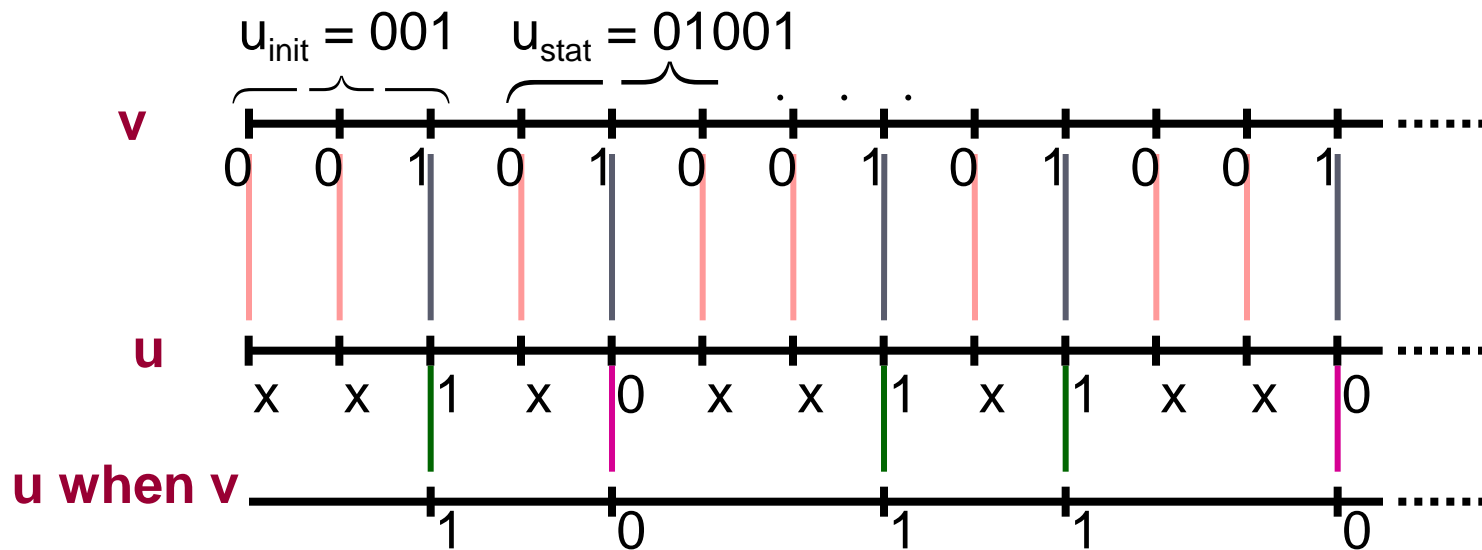


# when effects

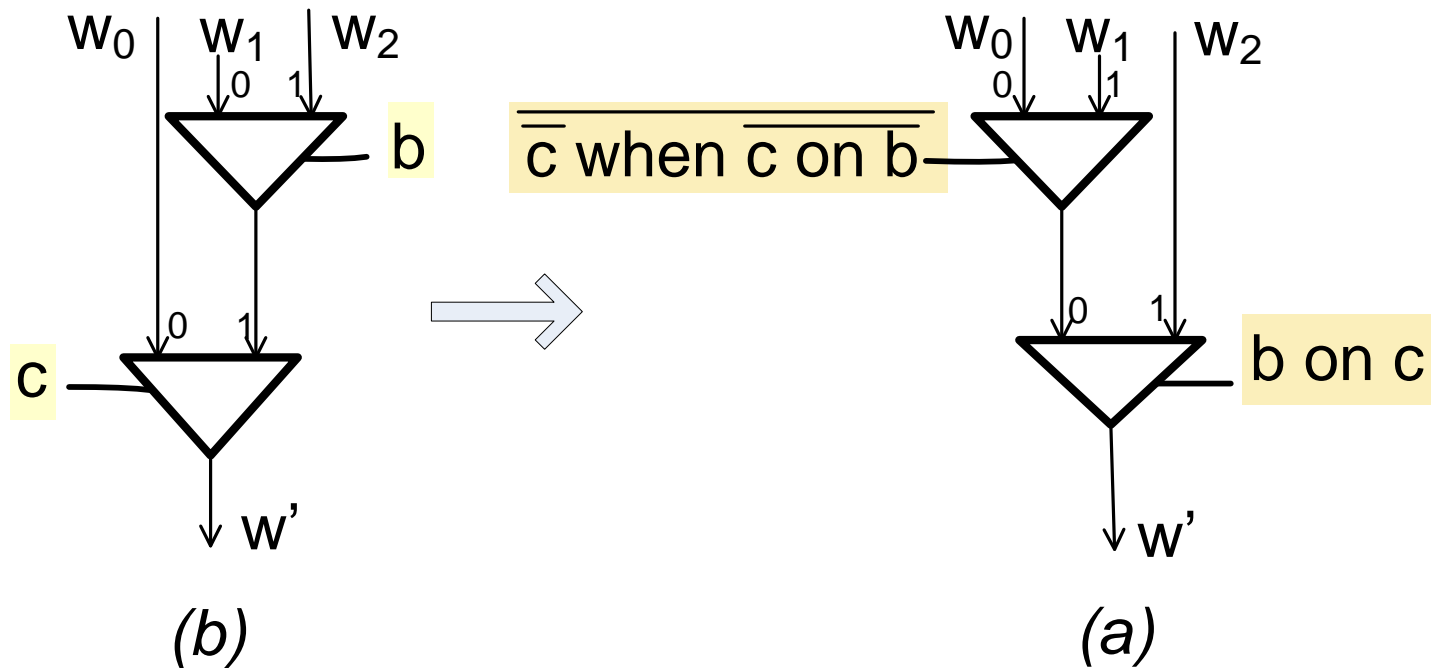
most meaningful whenever  
u subclock of v

$$(x.u) \text{ when } (0.v) = u \text{ when } v$$

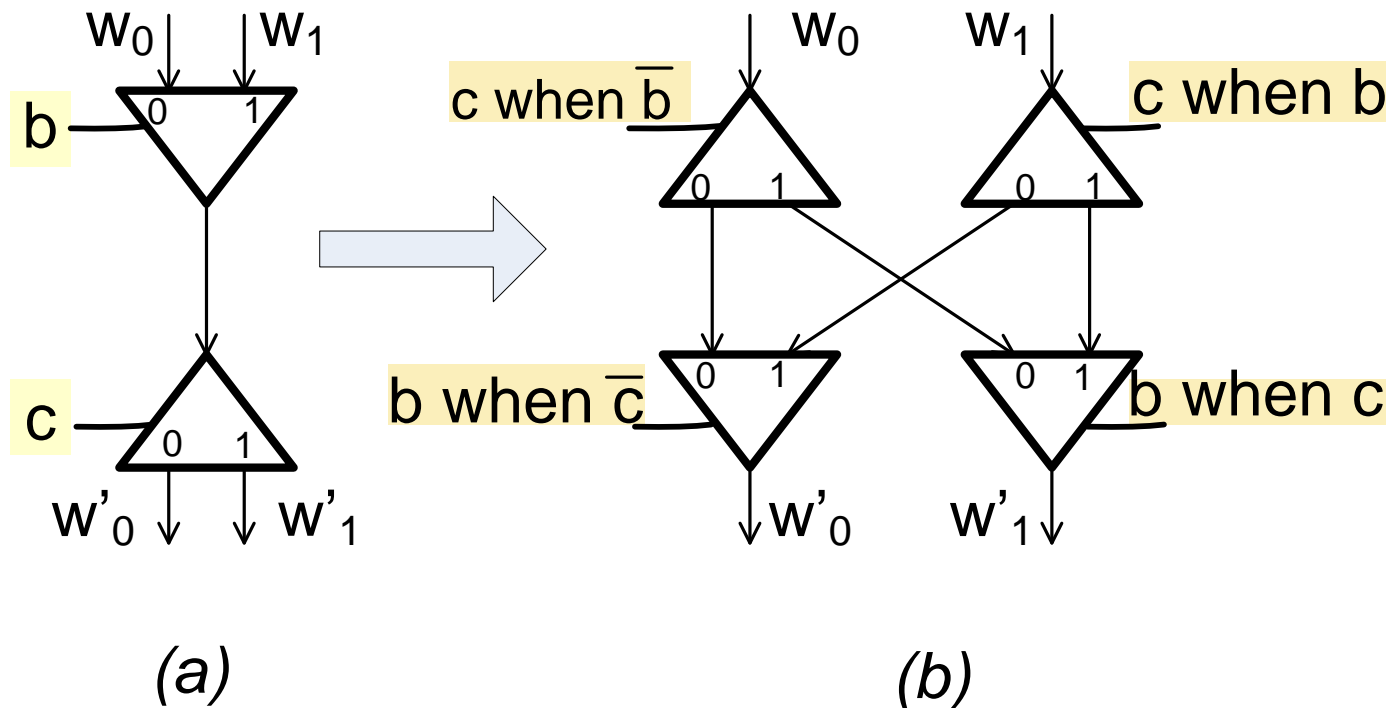
$$(x.u) \text{ when } (1.v) = x. (u \text{ when } v)$$



# Transitivity of Merges

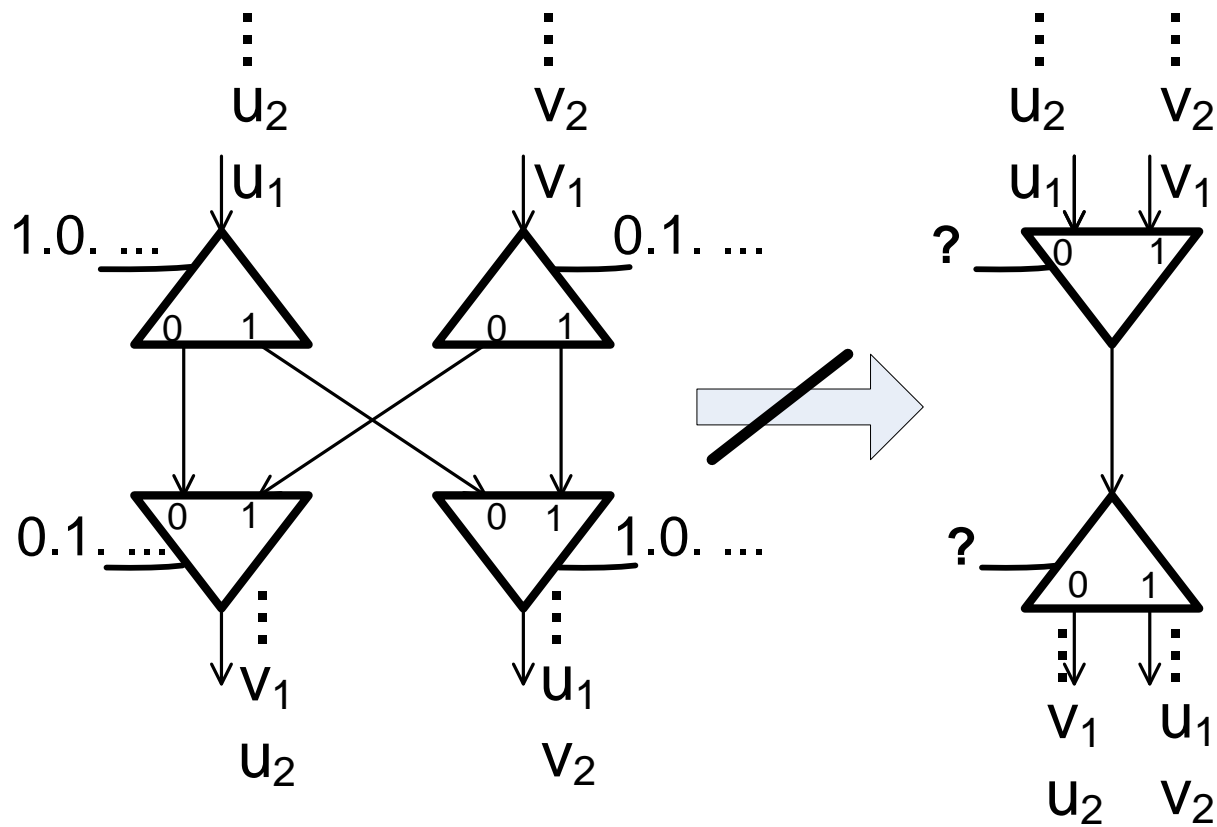


# Selects up across Merges

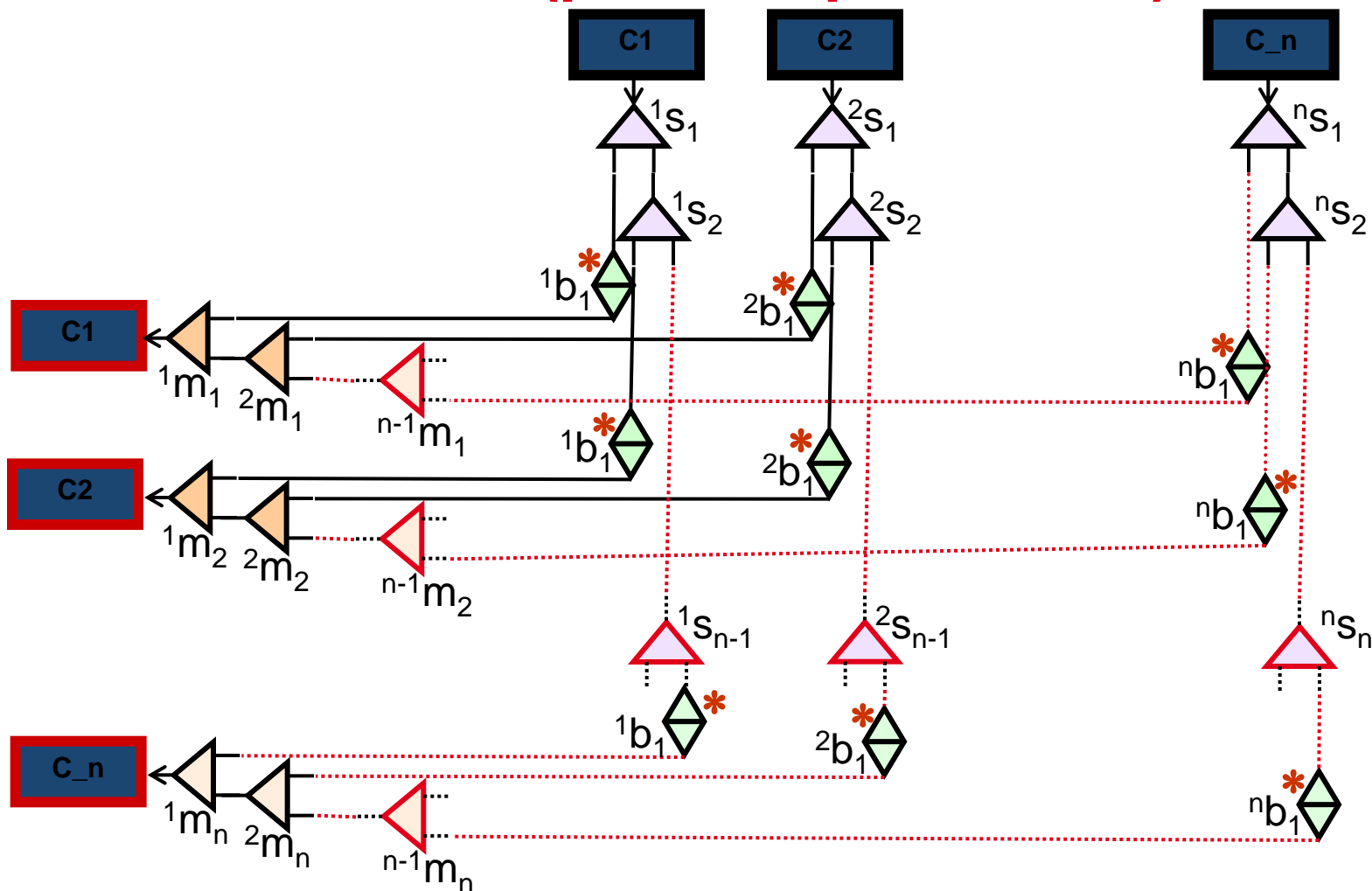




# Sharing vs disorder



# Normal forms (point-to-point links)



Regular scheduling arises naturally as solution space for optimal results

- Classical scheduling of Process Networks
- Scheduling of nested For-loops programs with affine bounds
- ...

Regular switching patterns can match the expressive level of description

- represent the transformations in data transfer/communication pattern
- play with the boundaries between data- and task- parallelism

# Task streaming level: Scheduling Theory

Original problem is to schedule independent tasks with hard deadline requirements on a mono- or multi-processor

- static scheduling of periodic tasks: Rate-Monotonic Analysis, Deadline-Monotonic
- dynamic scheduling of periodic/sporadic preemptible tasks: Earliest-deadline-first, least-laxity first,...

Simplified assumptions lead to positive methods/results

Then accounting for communications, context-switches in preemption, and further issues make life harder (and literature bigger)

First step is to specify in clear mathematical formulas what are the real-time constraints

# 5

## **A Clock Constraint Specification Language** Declarative specifications for the top floor

# Multiform Logical Time

Main idea is simple:

**every (pure) event that occurs repeatedly can be used as a *logical clock***

logical clock = sequence/flow of ticks/signals = activation condition

examples:

- (physical) clock cycles in a modern embedded processor
- ignition in a car (4 times each turn of engine, no matter the speed)
- sensor event detection (such as infra-red cells, gesture detection, lidars,...)
- *by the way, regular clockwatch physical time can also be considered logical...*

We claim Multiform Logical Time is natural and invaluable at design/specification time

- Expanding multiform logical time to uniform physical time depends on implementation conditions

# Clock Constraint Specification Language) CCSL

- Meant to express constraints and properties in (multiform) logical time
- Targets the platform-based design/AAA framework
- Inspired from Synchronous (actually polychronous) languages
- Inspired from Classical and Real-Time Scheduling notions.
- Not usually stand-alone (extracts the ordering relations between events whose meaning is part of a larger specification)
- Formal syntax to reason about logical time relations (including simultaneity)
- Concerns for expressivity and decidability

# Two natural partial orders

a subclock b ( $a \subseteq b$ )

- inspired from synchronous (polychronous) languages
- inspired from shapes of regular static scheduling / parallel allocation
- inspired from hardware and system design

a faster than b ( $a \leq b$ )

- inspired from real-time scheduling
- inspired from Timed Automata
- ...

used either in an imperative or declarative fashion  
each tractable individually, but only the combination truly expressive  
... and problematic



## Brief recap on Synchronous languages 1. Signal /Lustre

$x = \text{fct}(y, z)$

$x = \text{init} \rightarrow \text{pre}(y)$

$x = y \text{ default } z$

$x = y \text{ when } z$  ( $z$ : Boolean flow)

$c = \text{when } z, x = y @ c$

$c\_y = c\_z \quad c := c\_y$

$c\_x := c\_y$

$c\_x = c\_y \text{ union } c\_z$

$c = c\_z \text{ filtered\_by } ?z$  ( $?z$  Boolean sequence)

$c\_x = c\_y \text{ inter } c$

$\wedge x = \wedge y$

$\wedge x \# \wedge y$

$c\_x = c\_y$

$c\_x \# c\_y$

Logical clocks not syntactic elements in Lustre, no totally first-class in Signal, explicit activation condition in Scade

## Brief recap on Synchronous languages 2. Esterel

pause (next instant)

P;Q

P || Q

loop P end

present S then P else Q

if ?S then P else Q

emit S

signal S in P end

+ sequential computations on data variables

start P when S

abort P when S

stop P when S

suspend P when S

activate P when S

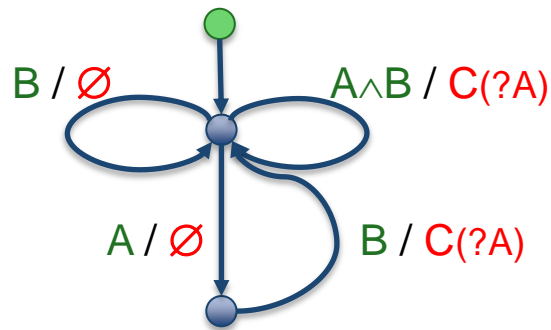
Signals are first-class citizen logical clocks

They represent shared variables with precise constructive consensus

semantics

## An example program

```
loop
  await A;
  await immediate B
  emit C(A?)
end
```

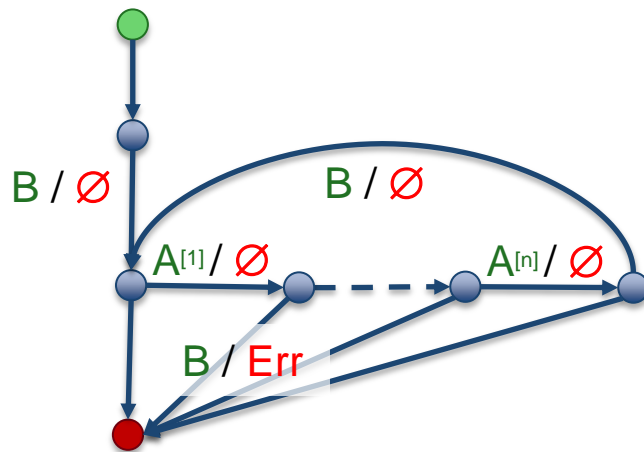


$C = A \text{ sampledOn } B$

- Logical clocks are the event-driven control structures (with registers that are latching instants)
- new clocks from old clock (only present-else case is problem, but usually triggered by some previous clock/register)
- Simple Mealy machine interpretation for each construct
- “Sensible” clocks should tick infinitely often

## A second example

```
(weak)every B
do (weak)abort
    await B; emit Err
when n*A
done
end
```



B SporadicOn A

- Constraints seen as Observers
- Assume/Guarantee

a, b, c .. clocks, w mask (infinite regular binary word)

## expressions

**a union b** only way to « speed up » rate (more ticks)

**a inter b** (less ticks, may become finite)

**a minus b** (less ticks, may become finite)

**fastest(a,b)** min of timings

**slowest(a,b)** max of timings

**sync(a,b)** when timings match

**a filteredBy w** subsampling, solution-oriented (k-periodic, regular)

**a sampledOn b** sporadic sampling

## relations/constraints

**a = b**

**a ≠ b**

**a ≤ b**

**a < b**

**a # b**

**a AlternatesWith b**

## Results (extracts)

Each individual constructs translates into a:

transition-labelled extended Büchi Mealy transition systems.

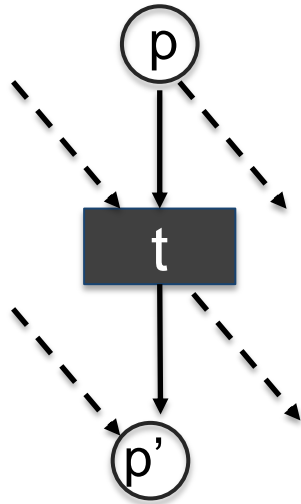
Not always Finite-State machine, also integer counters  
(cf  $\text{fastest}(C,D)$ , or  $C \leq D$  )

Full system is a parallel product of such t-E Büchi Mealy machines

Too Expressive (can encode Petri Nets with inhibitor arcs, and similar arguments as classical 2 counter machines)  $\rightarrow$  Turing-Complete

But many sources of sensible restrictions ensuring decidability

# Encoding asynchronous Petri Nets



## Clocks

- $t_p$ ,  $p_t$  for each  $t, p$  connected
- $t$  for each  $t$  (firing the transition)
- $p_{in}$  and  $p_{out}$  for put/get token to/from that place

$$t = \text{inter}(p_t, p \text{ in } \bullet t) \cap \text{inter}(t_p, p \text{ in } t \bullet)$$

$$t_1 \# t_2 \dots \# t_n$$

(transition are exclusive in time, not to mix tokens)

$$p_{in} = \text{union}(t_p, p \text{ in } \bullet t)$$

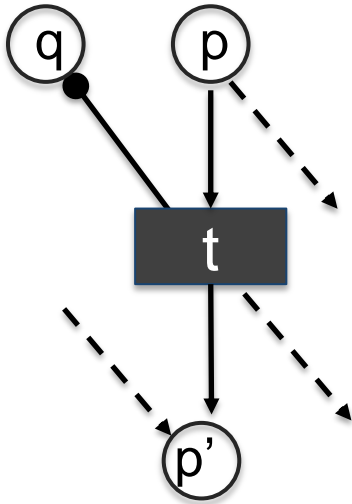
$$p_{out} = \text{union}(p_t, p \text{ in } t \bullet)$$

$$(p_{out} \text{ shiftedBy } \text{init\_Tok\_p}) \geq p_{in}$$

(tokens consumed only after being produced)

Note: in general Petri Nets, weights can be replaced by multiple places)

# Encoding PN inhibitor arcs



$$t = \text{inter}(\text{p\_t}, p \text{ in } \bullet t) \cap \text{inter}(\text{t\_p}, p \text{ in } t \bullet)$$

$$t_1 \# t_2 \dots \# t_n$$

*(transition are exclusive in time, not to mix tokens)*

$$\text{p\_in} = \text{union}(\text{t\_p}, p \text{ in } \bullet t)$$

$$\text{p\_out} = \text{union}(\text{p\_t}, p \text{ in } t \bullet)$$

$$(\text{p\_out shiftedBy init\_Tok\_p}) \geq \text{p\_in}$$

$$\text{p\_in} \# \text{p\_out}$$

$$\text{void\_q} = \text{fastest}((\text{q\_out shiftedBy init\_Tok\_q}), \text{q\_in}) \text{ inter } \text{q\_out}$$

*ticks when place becomes empty*

$$\text{unvoid\_q} = \text{void\_q} \text{ SampledBy } \text{q\_in}$$

*ticks when place becomes occupied*

$$t \text{ SampledOn } (\text{void\_q} \text{ union } \text{unvoid\_q}) = t \text{ SampledOn } \text{unvoid\_q}$$

*void alternates with unvoid, these main equation states that t occurs after the void and before the unvoid*



# Schedulability issue: simplest example

- $B \subseteq A$      $A$  filtered by (01) =  $B$  filtered by (01)

Here in fact  $B=A$ , since all even occurrences must coincide (while odd occurrences seem free),  $B$  can never catch up with any delay.

But it may be noticed too late in plain simulation (without backtracking)

# Many tractable subsets

If the clock dependency is a forest (no reconvergence)  
If it is a DAG but without synchronizing operators (inter,..)  
in the middle  
More to be defined  
...

# 6

## Conclusions



# Thank you



[www.inria.fr](http://www.inria.fr)

# 6

## Platform-Based Design and Y-Chart methodology

Sous-titre facultatif

