Technion

# The Plural Architecture
## Shared Memory Many-core with Hardware Scheduling

Ran Ginosar

Technion—Israel Institute of Technology

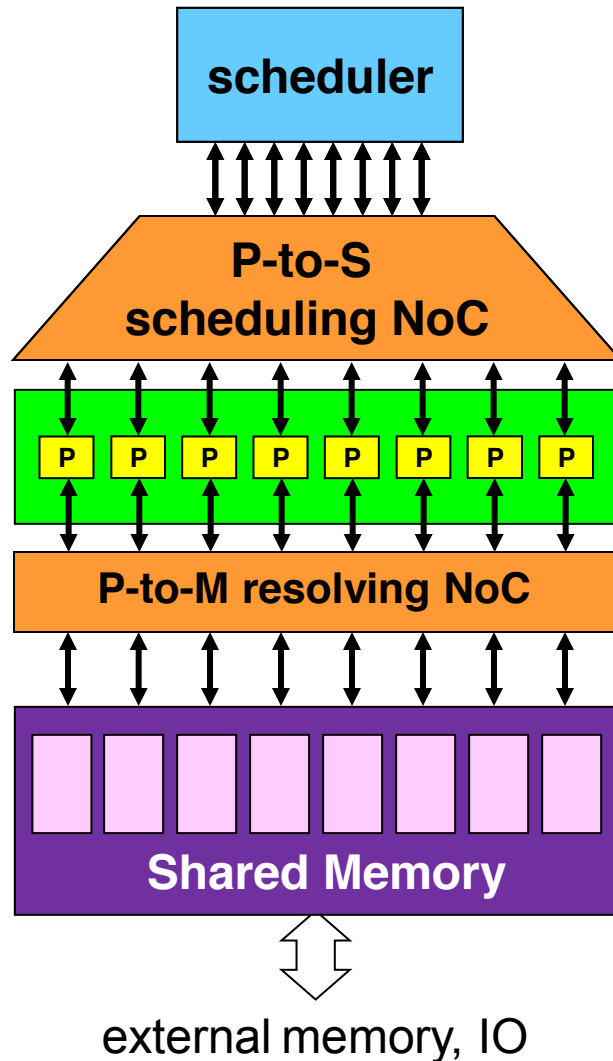&   Ramon Chips, Ltd.

Israel

MPSoC  July  2014

# The Plural Architecture



scheduler

P-to-S scheduling NoC

P P P P P P P P

P-to-M resolving NoC

Shared Memory

external memory, IO

Hardware scheduler / dispatcher / synchronizer

Low (zero) latency parallel scheduling enables fine granularity

Many small processor cores
Small private memories (stack, L1)

Fast NOC to memory
(Multistage Interconnection Network)
NOC resolves conflicts

SHARED memory, many banks
~Equi-distant from cores (2-3 cycles)

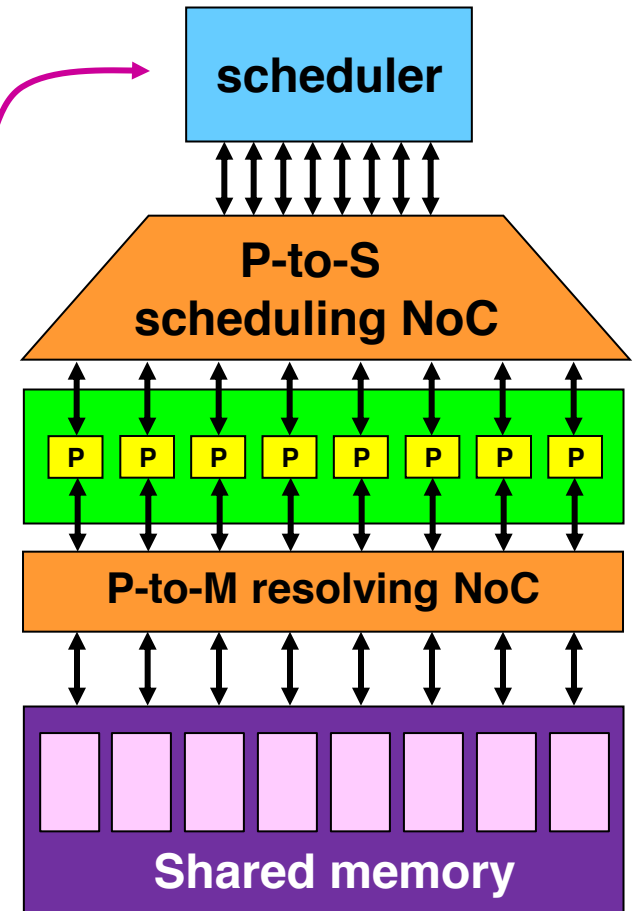"Anti-local" address interleaving
Negligible conflicts

# The Plural task-oriented programming model

- Programmer generates TWO parts:
  - Task-dependency-graph = 'task map'
  - Sequential task codes
- Task maps loaded into scheduler
- Tasks loaded into memory

**Task template:**

$$
\left.\begin{matrix} \textbf{singular} \\ \textbf{duplicable} \\ \textbf{control} \end{matrix}\right\} \quad \textbf{task } xxx\textbf{( dependencies )}
$$

```
{
        … # ….   // # is instance number
        …..
}
```



scheduler

P-to-S
scheduling NoC

P P P P P P P P

P-to-M resolving NoC

Shared memory

# Fine Grain Parallelization

Convert (independent) loop iterations

```
for ( i=0; i<10000; i++ ) { a[i] = b[i]*c[i]; }
```
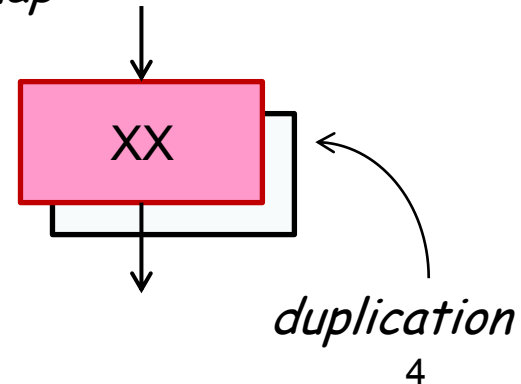
into parallel tasks

```
set_quota (XX,10000)

duplicable task XX
{ a[#] = b[#]*c[#]; }
 // # is instance number
```



Task map

XX

duplication

# Example: Linear Solver

# PIPELINED stream processing: **ManyFlow**

## Example: JPEG2000



**A** — serial

**B** — DWT (highly parallel *)

**C** — serial

**D** — Bit-plane encoding (highly parallel *)

**E** — serial

Num. cores utilized

Max 64 cores

Image compression time: 160 (relative time units)

Time

Low utilization: only 65%

# Hardware-like Pipeline

Time step i+1     Step i+2     Step i+3     Step i+4     Step i+5     Step i+6

Step i

Image k+4

Step i+7

Image k+5

Image k+1

Image k+6

Image k+2

Image k+7

Image k+3

Needs 5 stages: two with 64 cores each, three with one core each (total 131 cores)
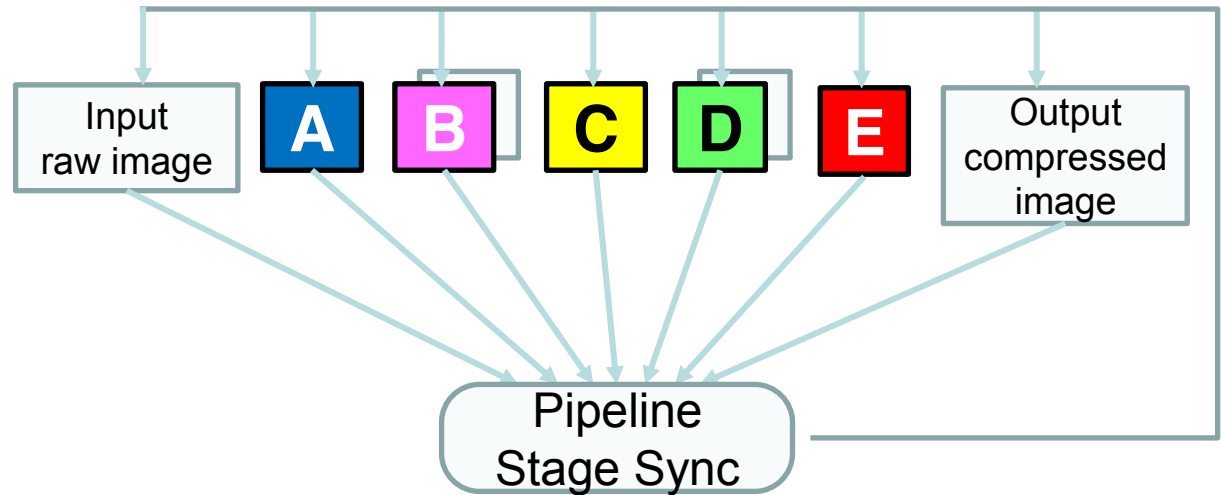If only 64 cores, time / step = 64x2 + 25 = 153 (how ? What is the utilization?)

Hard to program, inefficient, inflexible, fixed task per core. Need to store 5 images

# Parallel / pipelined ManyFlow

## Step i

Image k+4

Image k

Image k+1

Image k+2

Image k+3

All stages independent (order does not matter)
→ Can run concurrently
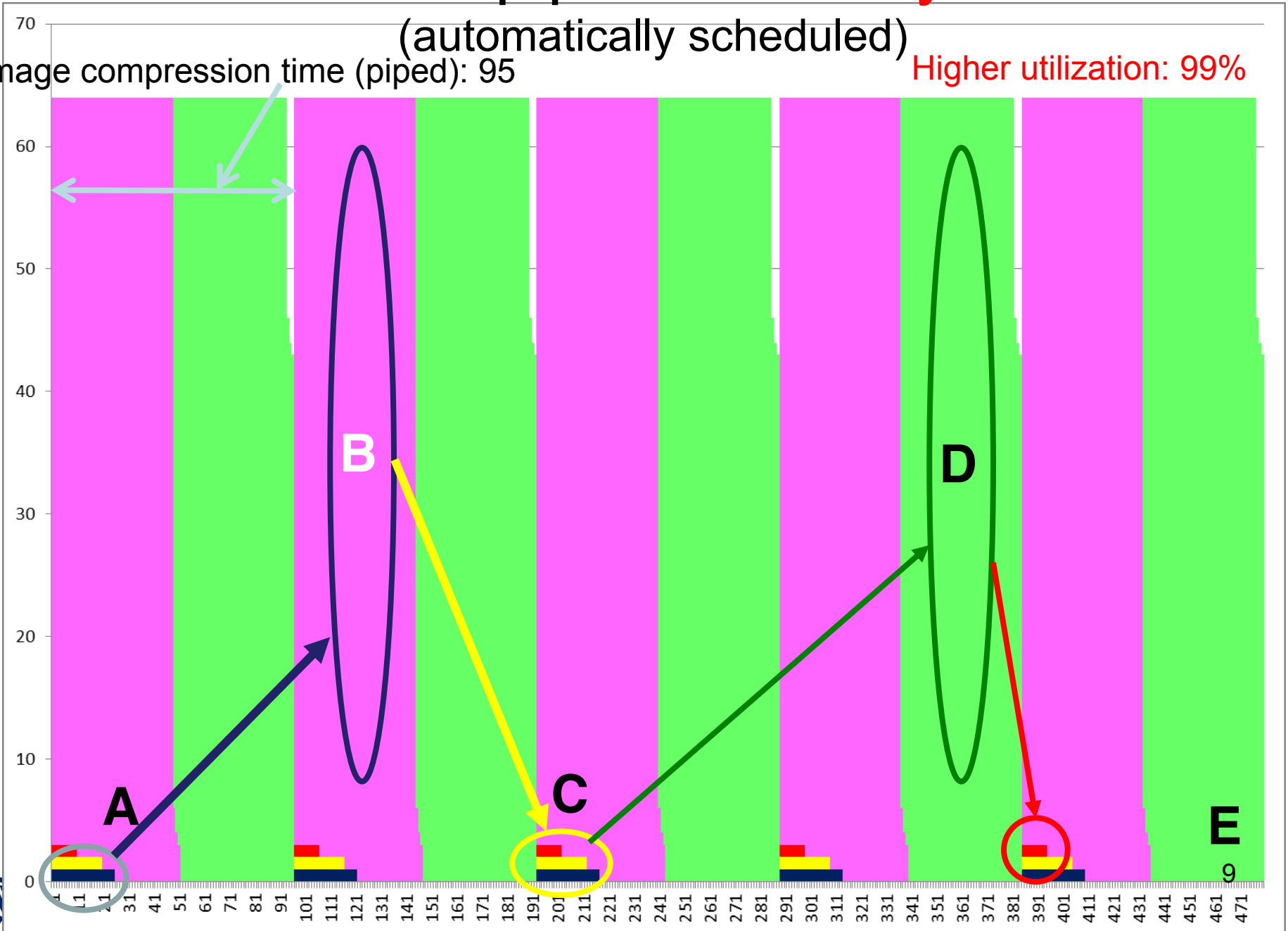→ Scheduler will dispatch most efficiently



| Input raw image | **A** | **B** | **C** | **D** | **E** | Output compressed image |

Pipeline Stage Sync

Bottleneck: need to store 7 images

# Parallel / pipelined ManyFlow
## (automatically scheduled)

Image compression time (piped): 95

Higher utilization: 99%



A    B    C    D    E

9

# The Plural Architecture: Some benefits

- Shared, uniform (~equi-distant) memory
  - no worry which core does what
  - no advantage to any core because it already holds the data
- Many-bank memory + fast P-to-M NoC
  - low latency
  - no bottleneck accessing shared memory
- Fast scheduling of tasks to free cores (many at once)
  - enables fine grain data parallelism
  - harder in other architectures due to:
    - task scheduling overhead
    - data locality
- Any core can do any task equally well on short notice
  - scales well
- Programming model:
  - PRAM-like
  - intuitive to programmers
  - "easy" for automatic parallelizing compiler & formal verification  (?)

# Summary

- Simple many-core architecture
  - Inspired by PRAM
- Hardware scheduling
- Task-based programming model
- Designed to achieve the goal of 'more cores, less power'