# Manycore Processor
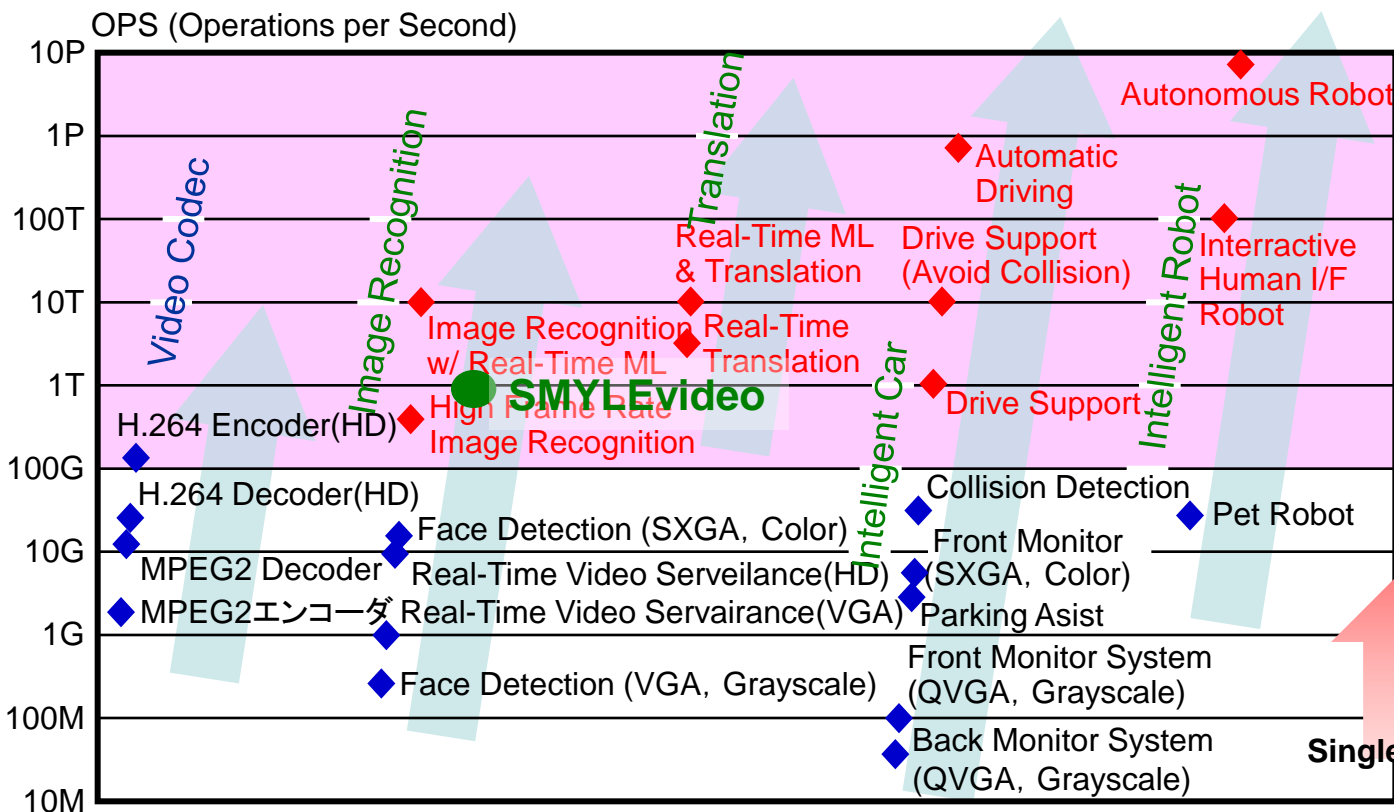# for Video Mining Applications



**Jan. 25th 2013**

Yukoh Matsumoto, Hiroyuki Uchida, Michiya Hagimoto,
Yasumori Hibi, Sunao Torii, Masamichi Izumida
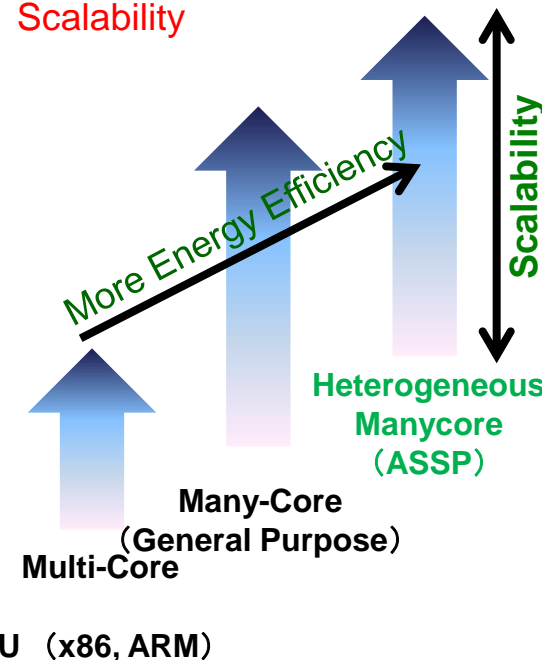**TOPS Systems Corp.**

# More Energy-Efficiency required for Next-Generation Embedded Systems

- **Next-Gen Embedded Systems** : **Requires more performance（100's GOPS~ 1 POPS）**
- **Power Consumption** ： **Already reached upper limit（～ W）**

**＜Expectations on Many-Core＞**

- 10's GOPS/W ⇒ 100's GOPS/W
- General Purpose ⇒ ASSP
- Scalability

OPS (Operations per Second)

| | |
|---|---|
| 10P | |
| 1P | Autonomous Robot |
| 100T | Automatic Driving |
| 10T | Real-Time ML & Translation / Drive Support (Avoid Collision) / Interactive Human I/F Robot |
| 1T | Image Recognition w/ Real-Time ML / Real-Time Translation / SMYLEvideo / Drive Support |
| 100G | High Frame Rate Image Recognition / H.264 Encoder(HD) |
| 10G | H.264 Decoder(HD) / Face Detection (SXGA, Color) / Collision Detection / Pet Robot |
| 1G | MPEG2 Decoder / Real-Time Video Serveilance(HD) / Front Monitor (SXGA, Color) |
| 100M | MPEG2エンコーダ / Real-Time Video Servairance(VGA) / Parking Asist |
| 10M | Face Detection (VGA, Grayscale) / Front Monitor System (QVGA, Grayscale) / Back Monitor System (QVGA, Grayscale) |

Video Codec — Image Recognition — Translation — Intelligent Car — Intelligent Robot

More Energy Efficiency

Scalability

Heterogeneous Manycore （ASSP）

Many-Core （General Purpose）

Multi-Core

Single CPU （x86, ARM）

Ref： NDEO Technology Roadmap 2009, I-48p, Fig. 1-6

**Energy-Efficient Computing goes to Heterogeneous & Manycore**

## "**ML and 3D Object Recognition** on **Image Stream**"

### Computer Vision (CW) ＝**EYE** (sensor) + **CEREBRUM** (decision)

**Decision**
- Recognition
- Learning

**Sensor**

**CPU**

Extremal Search
Refinement
Feature extraction
Matching

**Heterogeneous Manycore**

·**Distributed Processing**
·**Stream Computing**

TOPS

10
1
0.1
0.01
0.001   1GHz

Hetero Multi / Manycore
Computer Vision
Limit of Multi-Core
Multi-Core
Single CPU

2000    2015    2030

time

## Why we need Hetero Manycore?

- ■ **Conceptually:**           **Machine Learning (ML) ≒ Functional Configuration**
- ■ **High Perf. Requirement:** **More than 1 TOPS**
- ■ **Inherent Parallelism：**     **More than 99% of processing**
- ■ **Several types of Proc.：**  **Huge cost with Hardwires Implementation**
  - ➢ **Resolutions：**     **VGA, XGA, SGGA, FHD, 2K, 4K, etc.**
  - ➢ **Algorithms：**      **SIFT, Optical Flow, Ransac, Viola & Jones, Model based Recognitions, etc.**
  - ➢ **Others:**           **Multi-Medium Streams (MPEG-2, MPEG-4, H.264, etc.)**

### Key requirement is High Performance with Flexibility

**Camera(s)**
**Video(MPEG, H.264, …)**

**Data Parallel**
  ~n streams, 2-D image blocks
**Task Parallel**
  Gaussian Filter, Integral Map, Motion Vectors, Motion Compensation, Entropy Decode, iDCT, de-blocking Filter, etc.

**Data Parallel**
  Image Stream by frames, Motion Streams (MVs)

20%

**Feature Extraction（Low Level）**

**Video Decoder + α**

**Stream**     **Motion Stream**     **Audio Stream**

・**face**
・**building**
・**road**
・**speech**
・**music**

**Shot boundary Detection**

**Keyword Detection（Mid Level）**

$x_2$        $x_3$        ・・・$x_n$

**Event Detection（High Level）**

80%

・**Frame**
・**Shot**
・**Scene**
・**Video Sequence**

・**Highlight**
・**Person**
・**Dog**

**Data Parallel**
  Image Streams, Motion Streams, Pixels, Image blocks, frames, shots
**Task Parallel**; **Recognition Algorithms**
  Viola Jone, SIFT, SVM, Matching, optical flow, etc

**Recognized Images,**
**Extracted Images**

**Data Parallel**
  Search target, event detection
**Task Parallel**
  Matching, Editing

**More than 99% of application can be parallelized**

# Parallelisms in Algorithms for Video Mining

| Application | Objective | Algorithm | Parallelism | |
|---|---|---|---|---|
| Video Analysis | Prediction of Motion Vector | Optical Flow | Line | |
| | Specific feature detection and extraction | SIFT[*1] | Data Partitioning | |
| | Detection of human, and tracking | Cascaded Haar Like | Pixel Level | |
| | Line detection for field separation | Huff | Line | |
| | Elimination of error from continuous frames | Ransac[*2] | Sample Data | |
| Human Search | Face detection from several angles | Vector Face Detection | Pixel Blocks | |
| | Extraction of features on faces | Model Based Face Detection | Task Level | |
| | Specific feature detection and extraction | SIFT | Data Partitioning | |
| Video Editing | Segment Extraction | Graph based Segmentation | Grid Level | |
| | Detection of Motion Vector | Block Matching | Line | |

*1 SIFT：Scale-Invariant Feature Transform

*2 Ransac：Random Sample Consensus

**Many type of pallarelisms are inherent in algorithms for Video Mining**

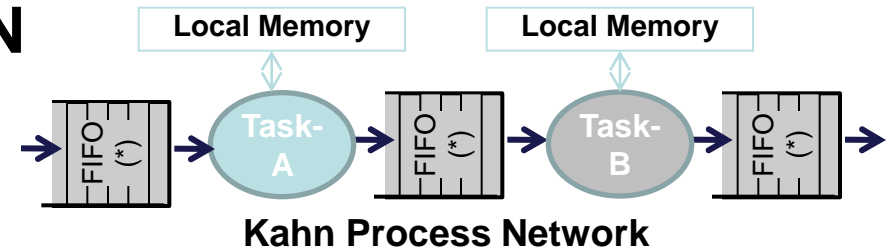# Goals of SMYLEvideo Manycare

- Real-Time Processing : **1TOPS~**

- Scalability : **10fps, 20fps, 30fps**

- Programmability : **Software Based Implementation**

- Flexibility : **OpenCV (Computer Vision)**

  SIFT, Optical Flow, Ransac, Viola & Jones, Model based Recognitions, SVM, etc.

- Low Power : **~1.5W**

- Low Clock Frequency : **~100MHz**

**High Performance, Programmable, Scalable**

**TOPS.**

■ **Distributed Processing** with KPN
 – **Non-Shared Memory Processes**
 – **Zero-Overhead Message Passing Mechanism**

**Local Memory** **Local Memory**

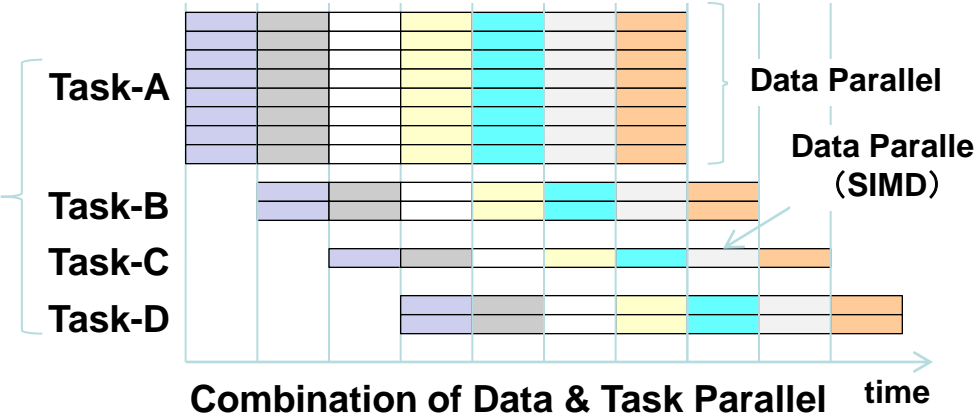FIFO (*) → Task-A → FIFO (*) → Task-B → FIFO (*) →

**Kahn Process Network**

■ **Combination of Parallelisms**
 – Distributed Parallel Processing（Task、Pipeline）
 – Data Parallelism（High-Level、Instruction Level）

**Task-A**

**Data Parallel**

**Data Parallel（SIMD）**

**Task Parallel**

**Task-B**

**Task-C**

**Task-D**

■ **Stream Processing** (Core)
 – Kernel
 – Stream-In (Read Message)
 – Stream-Out (Write Message)

**Combination of Data & Task Parallel** **time**

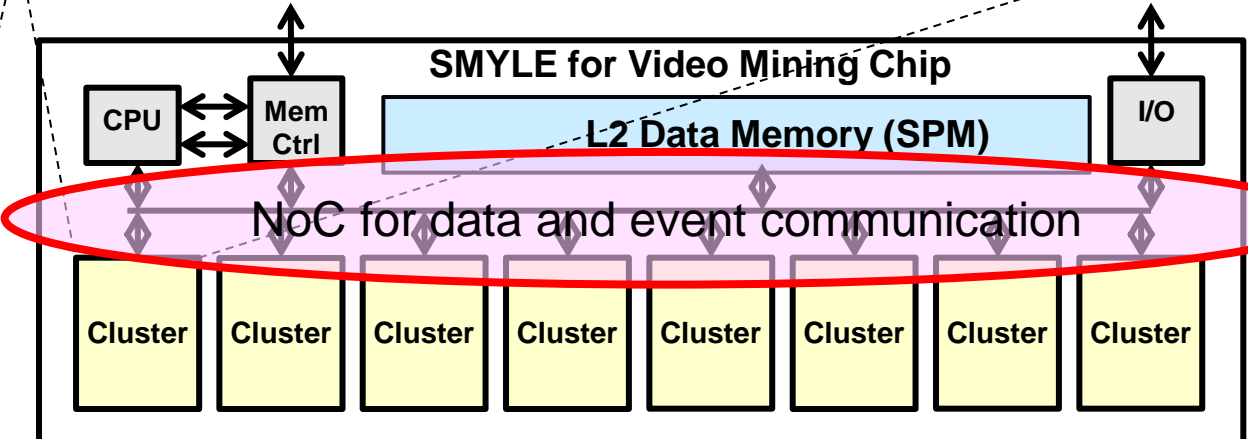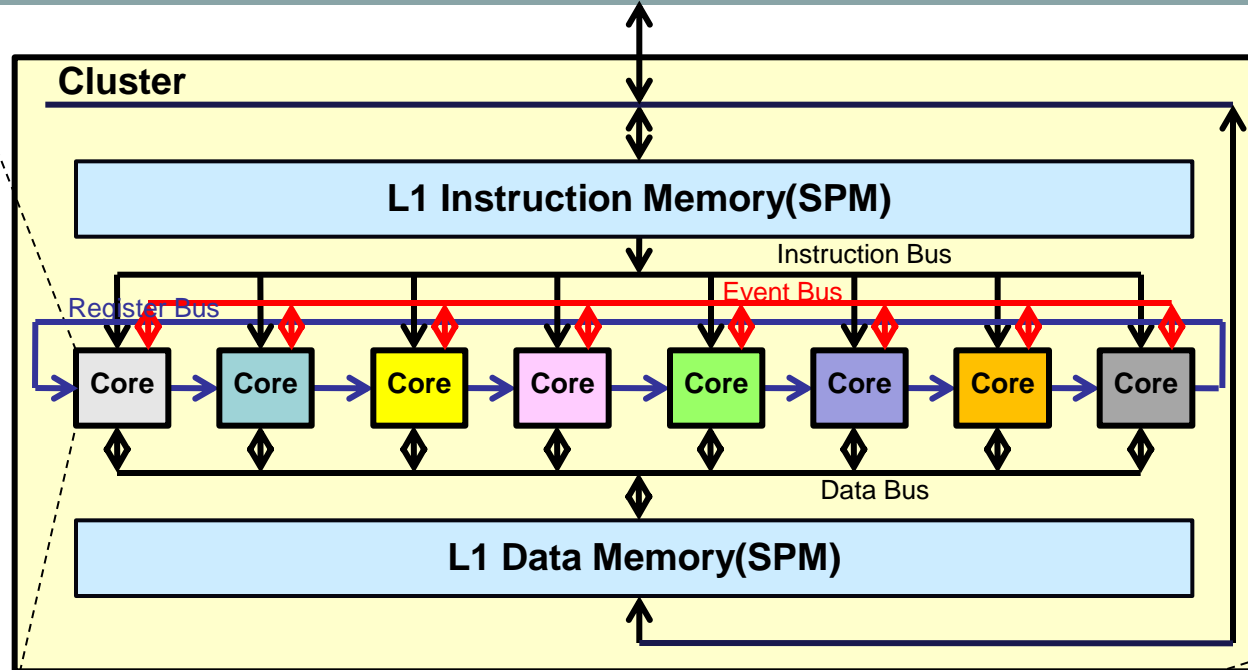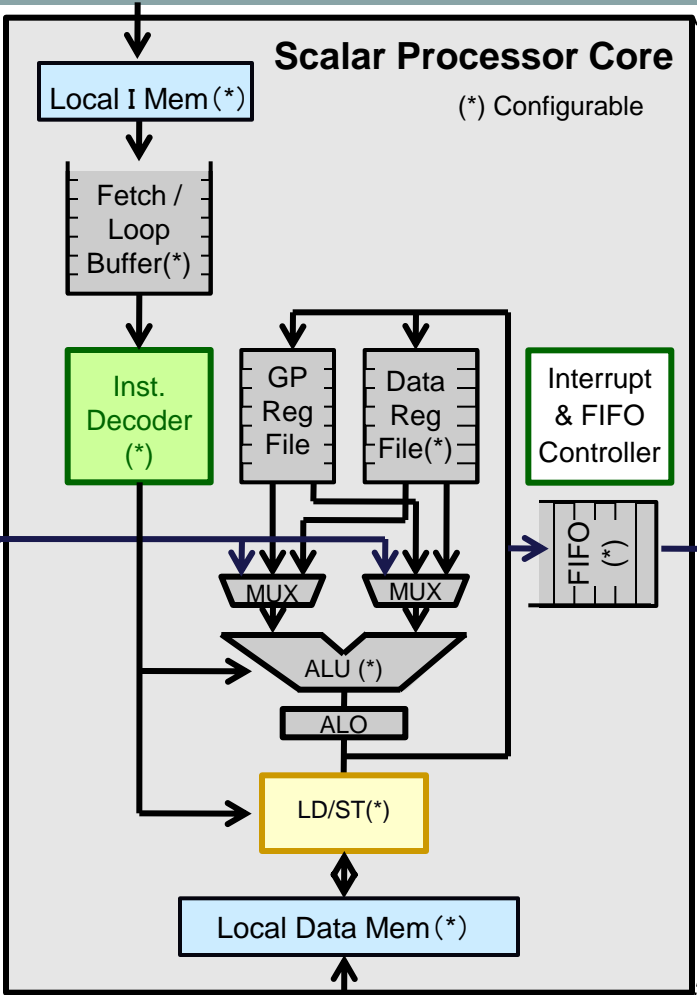| Stream IN | Stream IN | Stream IN | Stream IN | **Core can keep Processing of Kernel** |

■ **Optimization of Core**
 – Support Stream Processing : background Stream
 – Complex Inst : Reduction of Kernel cycle
 – FIFO support mechanism
 – Reduction of energy for instruction / data supply

| Kernel | Kernel | Kernel | Kernel |

Stream OUT | Stream OUT | Stream OUT | Stream OUT

**Distributed Processing, ZOMP, Task Parallel, Stream Processing , ASIP**

時間

# SMYLEvideo : Basic Architecture

## Scalar Processor Core

(*) Configurable

- Local I Mem(*)
- Fetch / Loop Buffer(*)
- Inst. Decoder (*)
- GP Reg File
- Data Reg File(*)
- Interrupt & FIFO Controller
- FIFO (*)
- MUX
- MUX
- ALU (*)
- ALO
- LD/ST(*)
- Local Data Mem(*)

> **Local Inst. Memory Configuration**
> **Loop Buffer Configuration**
> **Instruction Externtion**
> **(Decoder, ALU, LD/ST)**
> ~~Data Register Configuration~~
> 
> **Local Data Memory Configuration**

## Cluster

**L1 Instruction Memory(SPM)**

Instruction Bus

Event Bus

Register Bus

Core | Core | Core | Core | Core | Core | Core | Core

Data Bus

**L1 Data Memory(SPM)**

## SMYLE for Video Mining Chip

CPU | Mem Ctrl | **L2 Data Memory (SPM)** | I/O

NoC for data and event communication

Cluster | Cluster | Cluster | Cluster | Cluster | Cluster | Cluster | Cluster

## Application Domain Specific Scalable Heterogeneous Manycore
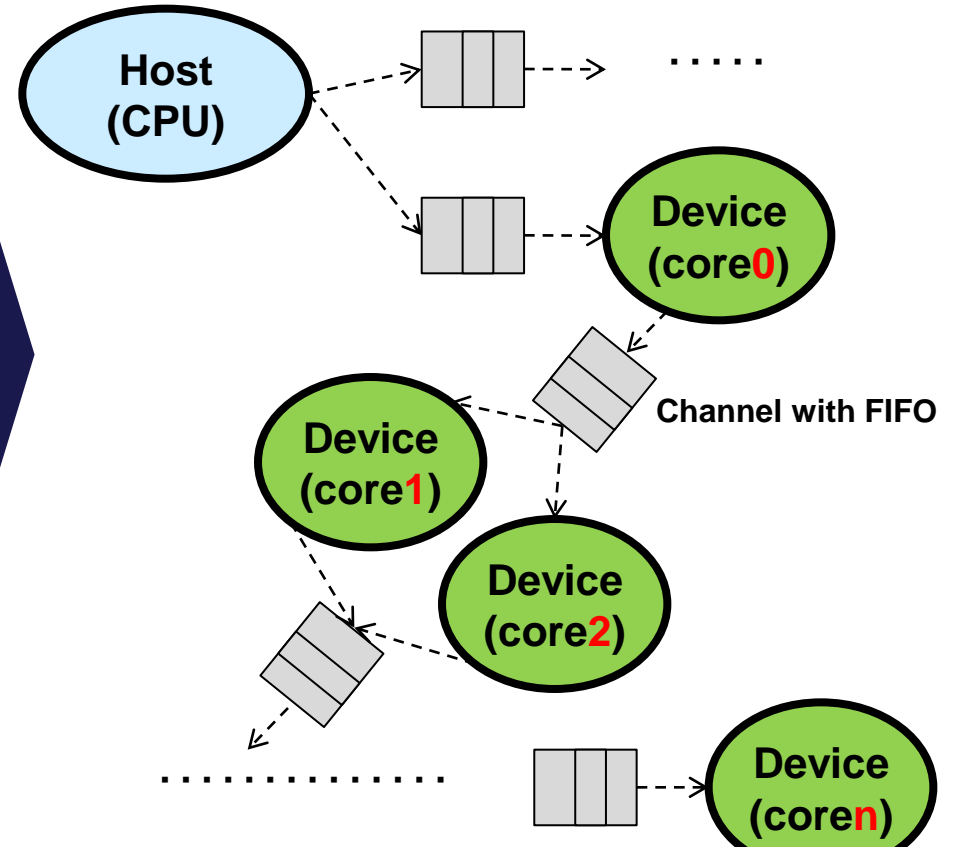
**OpenCL (CPU centric)**

- Bottleneck
  - Processing on Host
  - Increasing communication with Host
- **Hard to express distributed processing**
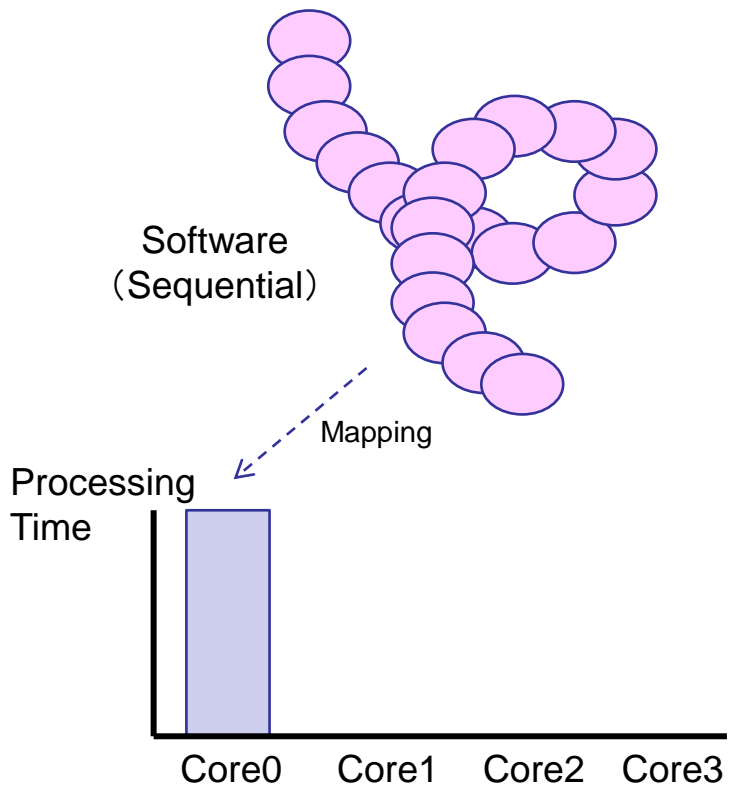
**Distributed Processing**

- **Scalability**
- **Can combine with Data Parallel Processing**



**Take an approach of Distributed Processing for removing bottlenecks**

Software
（Sequential）

Mapping

Processing
Time

Core0    Core1    Core2    Core3

**Distributed
Processing**

☐ Profiling
☐ Streaming
☐ Partitioning
☐ Grouping
☐ Optimizing

Software （Fine Grain Distributed Processing）

**Grouping**    **Adjusting to No. of
cores**

PE0
PE1
PE2

t

Mapping

Processing
Time

PE0
PE1
PE2

**Optimize**

Fine Grain is
better for load balancing

Core0    Core1    Core2    Core3

＜Issues＞

＜Advantages＞

- Cannot utilize **Multi-Core / Many-Core**
- **Requires High Performance** for Video Mining

- Utilize Multi-Core / Many-Core resources
- **Easy to balance the load**

**Investigation has done on Many Algorithms ; Viola & Jones, SVM, SIFT, etc.**

- **■** Goal :      **High-Performance & Low-Power Programmable Accelerator**

    (Energy-Efficient, Low Cost, Flexible, Scalable)

- **■** Approach :    **Low Clock Frequency**

◆ <u>Power consumption</u>

$$P_{Total} = P_{Dynmic} + P_{Static}$$

$$= \tfrac{1}{2}\, \alpha C\, V^2 f\, \alpha\ +\ I_{Leak}\, V$$

C: Load capacity
V: Source voltage
f: Frequency
α: Siwitching rate
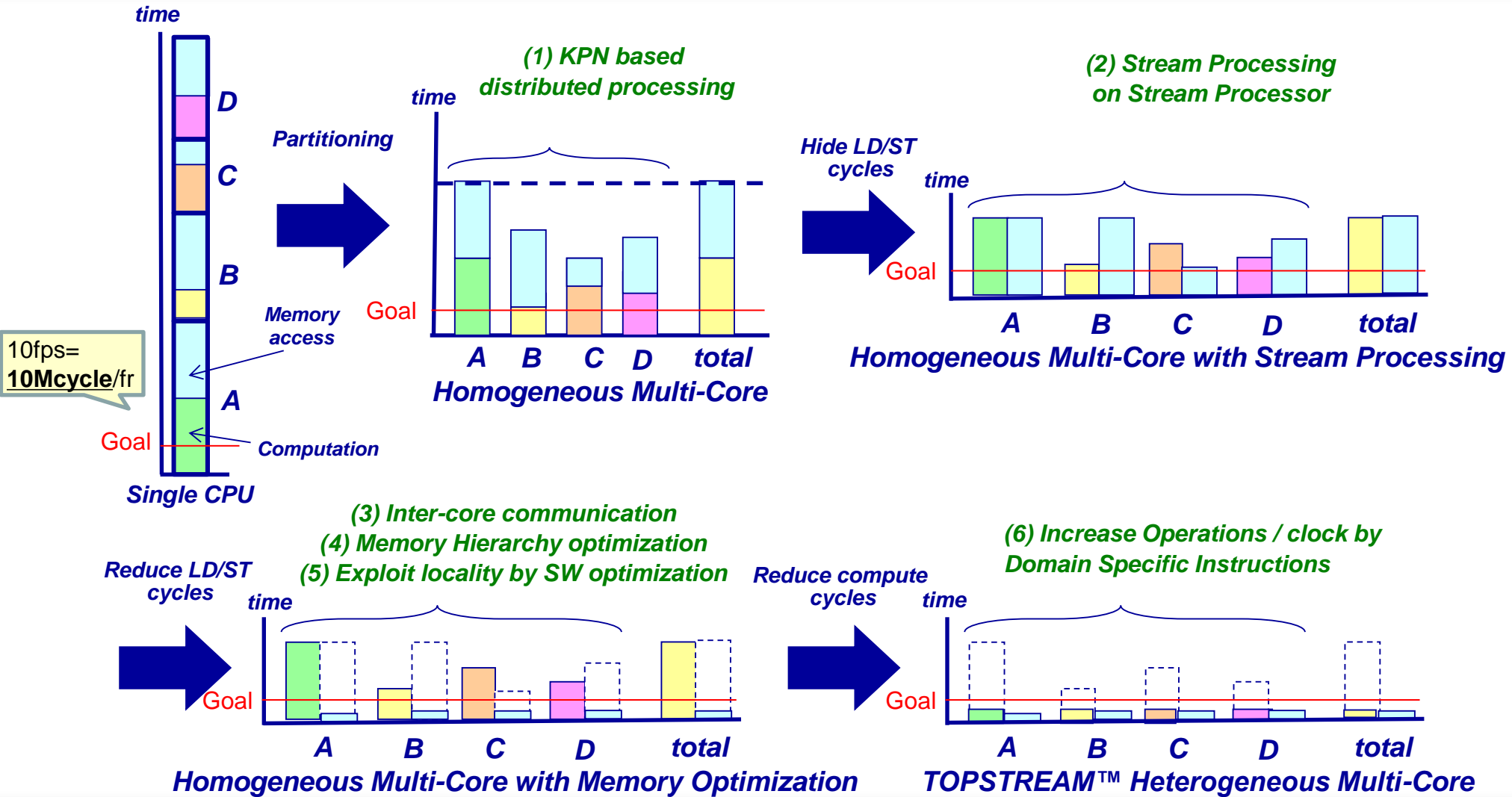$I_{Leak}$: Leakage current

◆ <u>Performance</u>

$$\text{Performance} = \text{OPC} \times f$$

OPC: Operation Per Clock

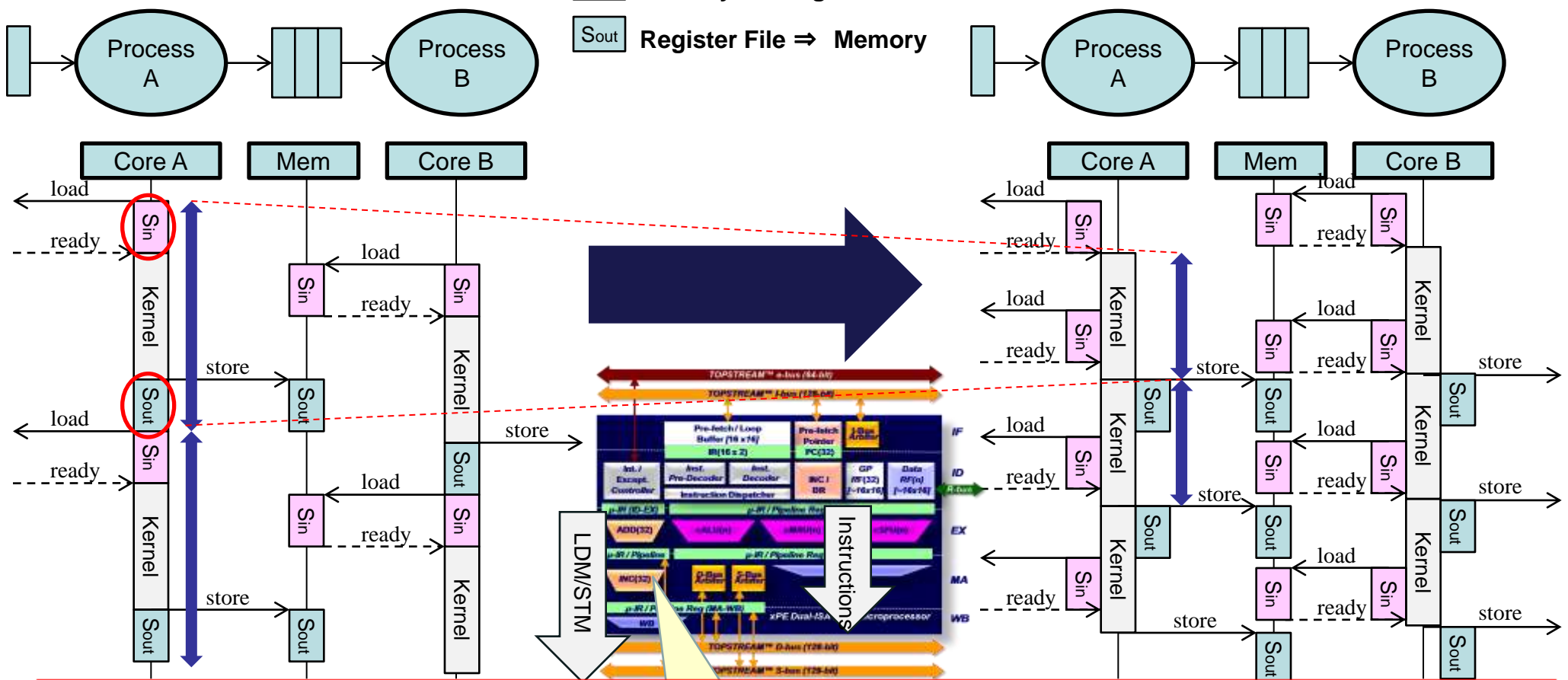**High Performance @ Low Clock Frequency drives Low Power**

# Approach to reduce clock frequency with Architecture-Algorithm Co-Design



time

D
C
B
A

Partitioning

Memory access

**10fps=**
**10Mcycle/fr**

Goal

Computation

**Single CPU**

*(1) KPN based distributed processing*

time

Goal

A    B    C    D    total
**Homogeneous Multi-Core**

Hide LD/ST cycles

*(2) Stream Processing on Stream Processor*

time

Goal

A    B    C    D    total
**Homogeneous Multi-Core with Stream Processing**

Reduce LD/ST cycles

*(3) Inter-core communication*
*(4) Memory Hierarchy optimization*
*(5) Exploit locality by SW optimization*

time

Goal

A    B    C    D    total
**Homogeneous Multi-Core with Memory Optimization**

Reduce compute cycles

*(6) Increase Operations / clock by Domain Specific Instructions*

time

Goal

A    B    C    D    total
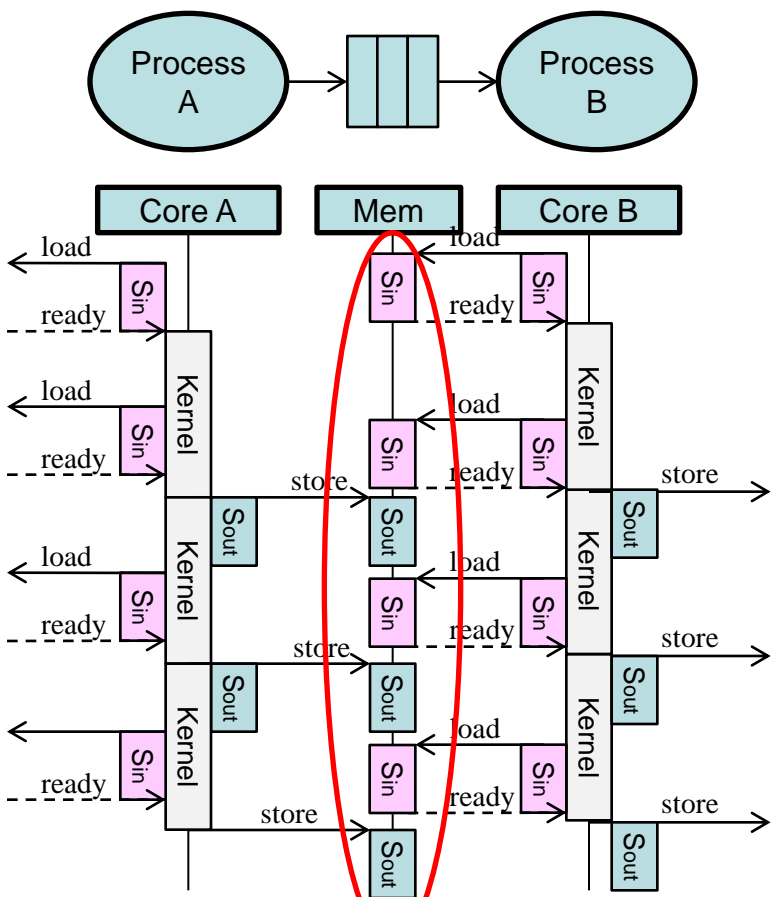**TOPSTREAM™ Heterogeneous Multi-Core**

■ Hide overhead of Stream-In and Stream-Out

$S_{in}$   Memory ⇒ Register File
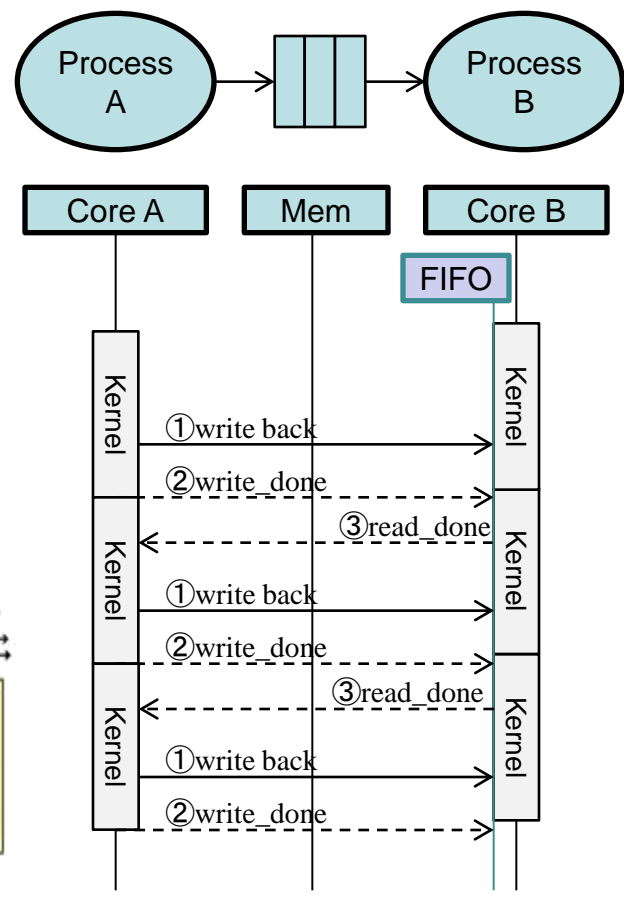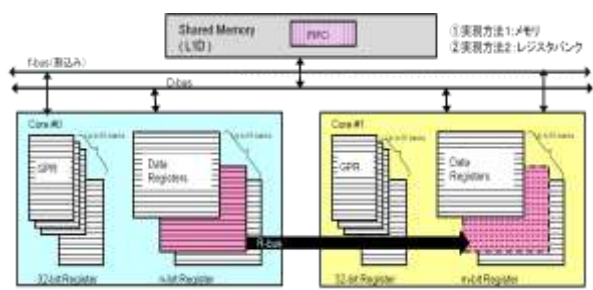
$S_{out}$   Register File ⇒ Memory



**Hide cycles for inter-process communication (Stream-In and Stream-Out)**

■ Reduction of Memory Access Bandwidth and its Energy



**Register Bank Sharing**

**Reduce memory access time and Energy!**

# Reduction of memory traffic

■ **Path for Message Passing**



Via Mem → Software to implement FIFO on memory

Via Reg → Software to implement FIFO on register

Total Memory access to Local Memory

- Via
- Via Reg

JPEG Decoder, JPEG Encoder, H.264 Decoder, Ray Tracing

**Significant Reduction of Memory Trafifc : more than 30%**

# Zero-Overhead Message Passing Mechanism

■ Remove cycles and memory access for checking FIFO counts and synchronization



**＜Hardware＞**
・FIFO Configuration Reg
・FIFO Counter

**FOWAIT0:FINT0:**STM [R1],D0x16

**＜Software＞**
・Prefix Instruction
　-FISYNCn: Check & Wait for Event
　-FOWAITn: Check & Wait for Outpur
　-FOINT:　　Event Generation
　-FIINT:　　Event Generation

**FISYNC0:FINT0:**LDM D0x16, [R0]

**No overhead : Just Add Prefix instructions on Stream-In & Stream-Out**

# Memory Access reduction by Distributed Stream Processing

- **Memory Centric Processing**
  - Each core works data on External Memory
  - Integration of processors and memories

- **Distributed Stream Processing**
  - **Core to Core Stream passing**
    - **On-Chip memory**
    - **Register Sharing**



Reduction of Memory Bandwidth Requirement

Ext. Mem

Core

External Memory

**Bottleneck**

Core 0    Core 1    Core 2    Core 3

**Frame based processing**

Eternal Memory

Core 0    Core 1    Core 2    Core 3

**Block based processing**

**Increase scalability by reducing memory bandwidth requirement**

# Frame based vs. Block based Processing

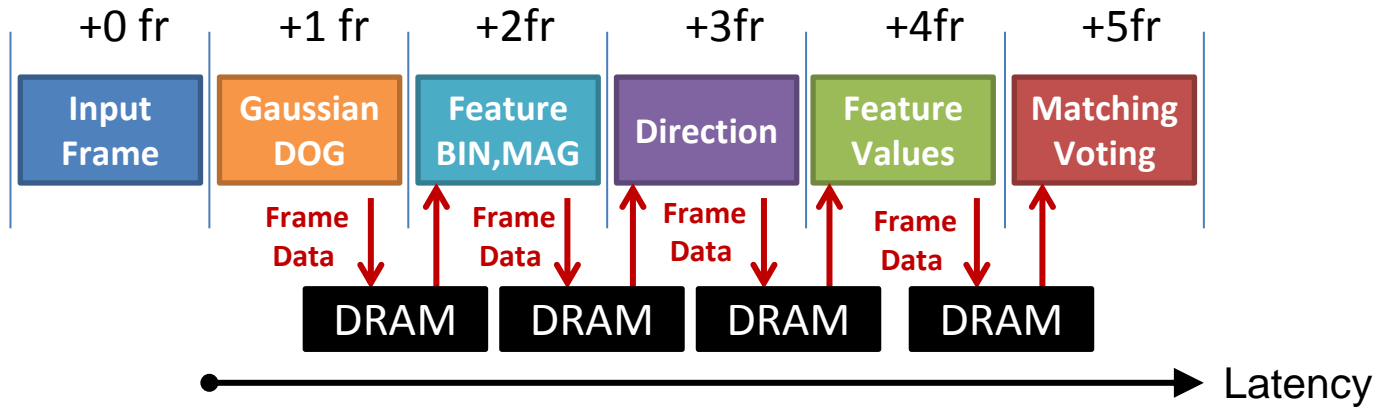Ex). SIFT

## Frame based processing

| +0 fr | +1 fr | +2fr | +3fr | +4fr | +5fr |
|---|---|---|---|---|---|
| **Input Frame** | **Gaussian DOG** | **Feature BIN,MAG** | **Direction** | **Feature Values** | **Matching Voting** |

Frame Data · Frame Data · Frame Data · Frame Data

DRAM · DRAM · DRAM · DRAM

→ Latency

## Block based processing

| +0fr | +1fr | +2fr | +3fr |
|---|---|---|---|
| **Input Frame** | | Input Frame | |

**Gaussian +DOG** · Gaussian +DOG

Block Data — **Feature Extraction**

Block Data — **MAG,BIN Generation**

Block Data — **Direction**

Block Data — **Feature Values**

Block Data — **Matching Voting**

→ Latency
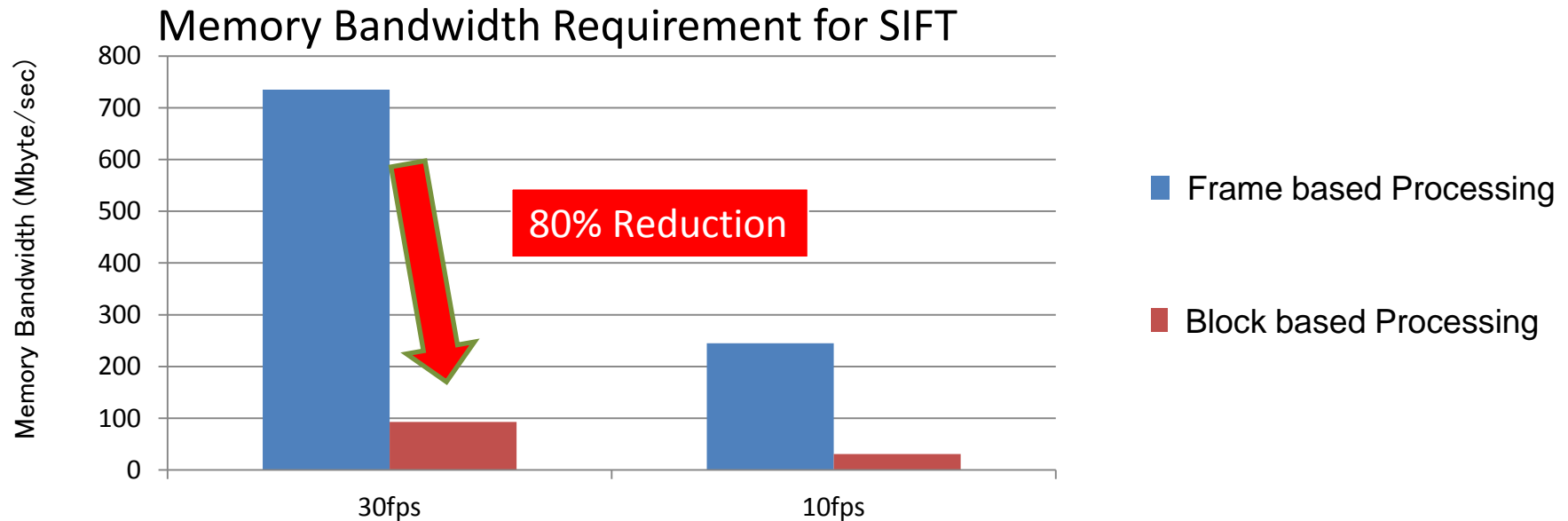
## Smaller latency with Smaller Data

# Frame based vs. Block based Processing

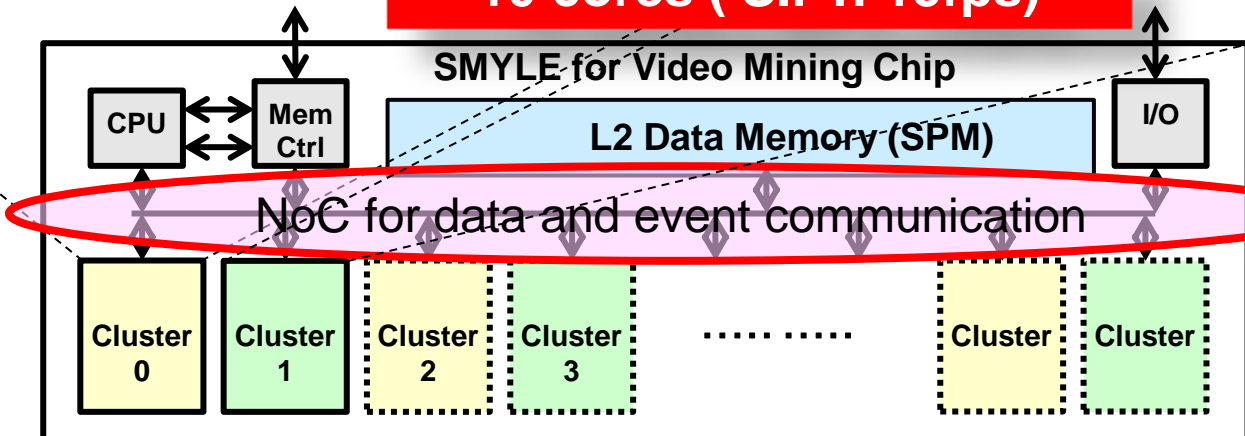| | Frame based Processing | Block based Processing |
|---|---|---|
| Global Memory Usage | 22Mbytes | 3.1Mbytes |
| Cluster Local Memory Usage | 0.15Mbytes | 0.8Mbytes |
| Recognition Latency | 167mSec(5.1frame) | 100mSec(3frame) |

## Memory Bandwidth Requirement for SIFT



**80% Reduction**

- Frame based Processing
- Block based Processing

**Memory Usage : 1/7,  Memory Bandwidth Requirement : 1/5**

# SMYLEvideo Configuration

**Cluster 0**

L1 Instruction Memory(SPM)

Instruction Bus

Event Bus

Register Bus

| QVP C0 | QVP C1 | QVP C2 | QVP C3 | SVP C4 | SVP C5 |

Data Bus

L1 Data Memory(SPM)

**Cluster 1**

L1 Instruction Memory(SPM)

Instruction Bus

Event Bus

Register Bus

| SVP C0 | SVP C1 | QVP C2 | QVP C3 |

Data Bus

L1 Data Memory(SPM)

**10 cores ( SIFT: 15fps)**

SMYLE for Video Mining Chip

CPU

Mem Ctrl

L2 Data Memory (SPM)

I/O

NoC for data and event communication

| Cluster 0 | Cluster 1 | Cluster 2 | Cluster 3 | ·········· | Cluster | Cluster |

QVP : Quad V processor (256-bit core)

SVP : Single V processor (64-bit core)

**Scalable Performance and Functionality with adding Clusters**

# Conclusions

- **Manycore** will play a crucial role in extending the roadmap for enabling the next generation **SoCs** required for **"Video Mining"** one of Computer Vision systems.

- **Zero-Overhead Message Passing Mechanism (ZOMP)** can **efficiently increases the system performance** and **scalability** of Manycore processors.

- **Block based distributed processing** drastically **reduces memory access bandwidth** and **increases room for higher performance** on Manycore processors.

- **SMYLEvideo** provides **scalability in performance and functionality** with its **clustered architecture**.