# Multi-Core and GPGPU Acceleration of Video Coding
## -- 100x Speedup of Motion Estimation --

**Wei-Chih Chen      Youn-Long Lin**

# National Tsing Hua University

# Outline

- **Introduction**
- **System and Software**
- **Proposed Methods and Implementations**
- **Experiment Results**
- **Conclusion**
- **Reference**

- **Introduction**
- **System and Software**
- **Proposed Methods and Implementations**
- **Experiment Results**
- **Conclusion**
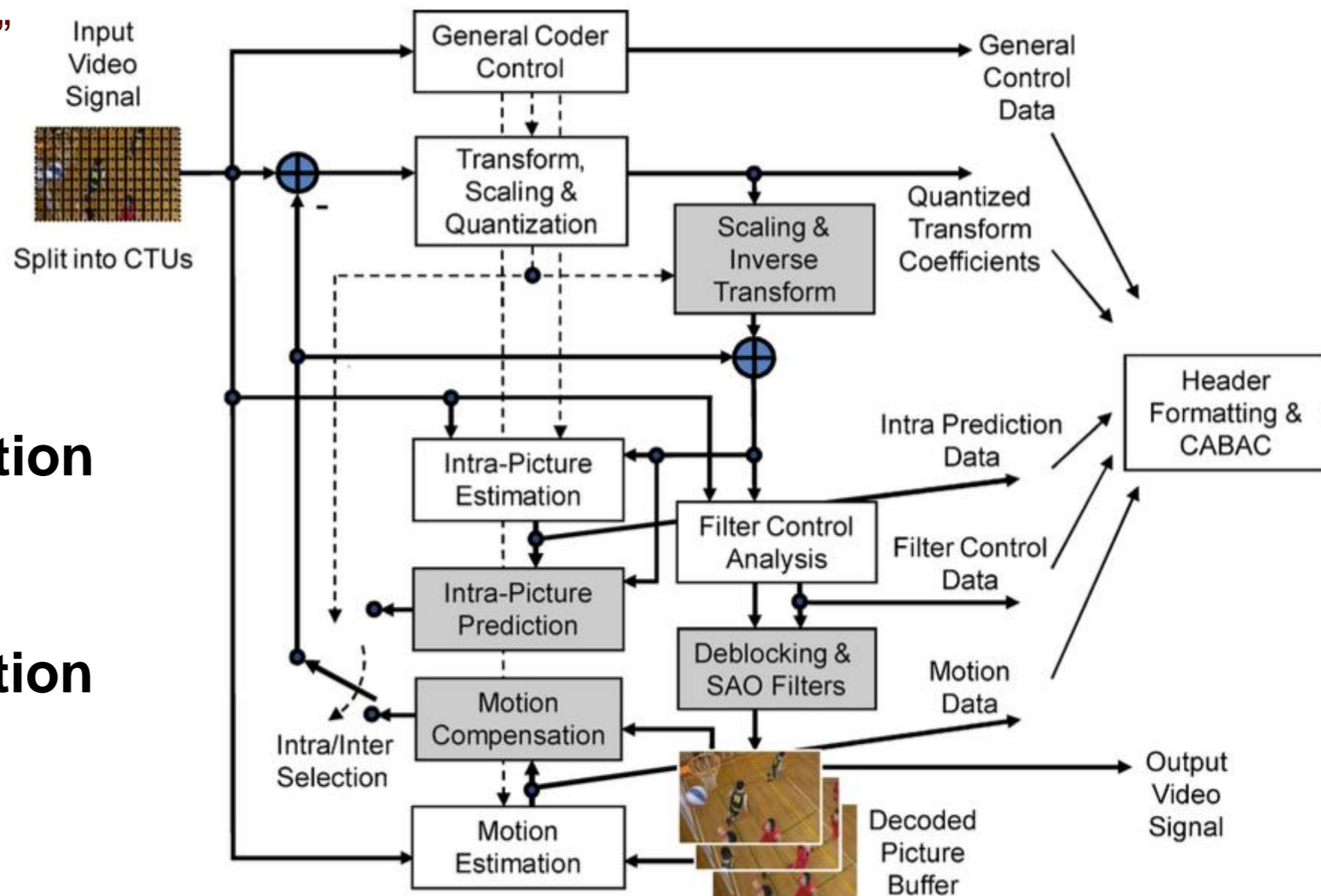- **Reference**

# Block-Based Hybrid Coding

- **Contemporary video coding**
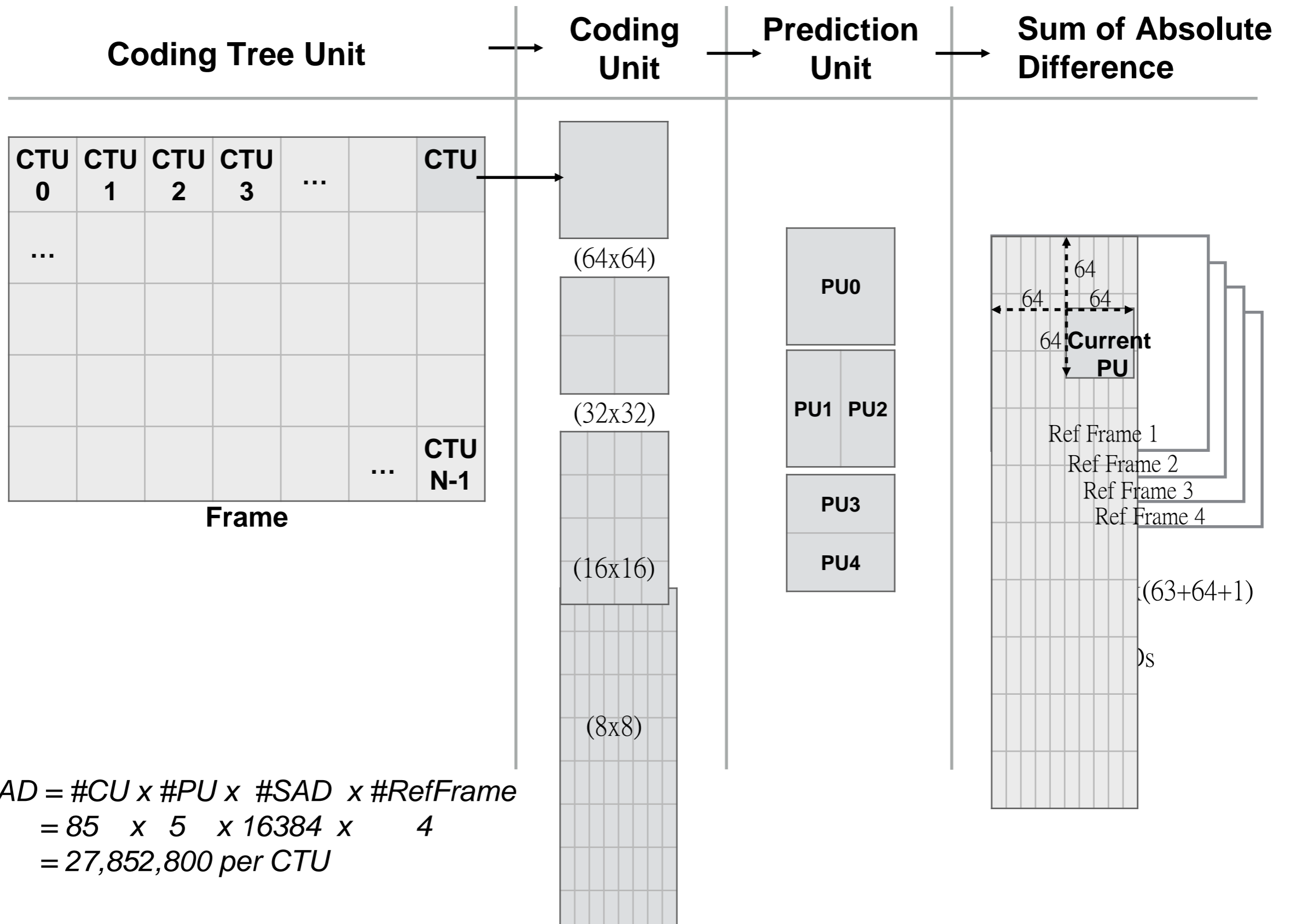  - Uses the same "hybrid" approach starting from H.261

- **Intra-picture prediction**
  - Spatial correlation

- **Inter-picture prediction**
  - Temporal correlation

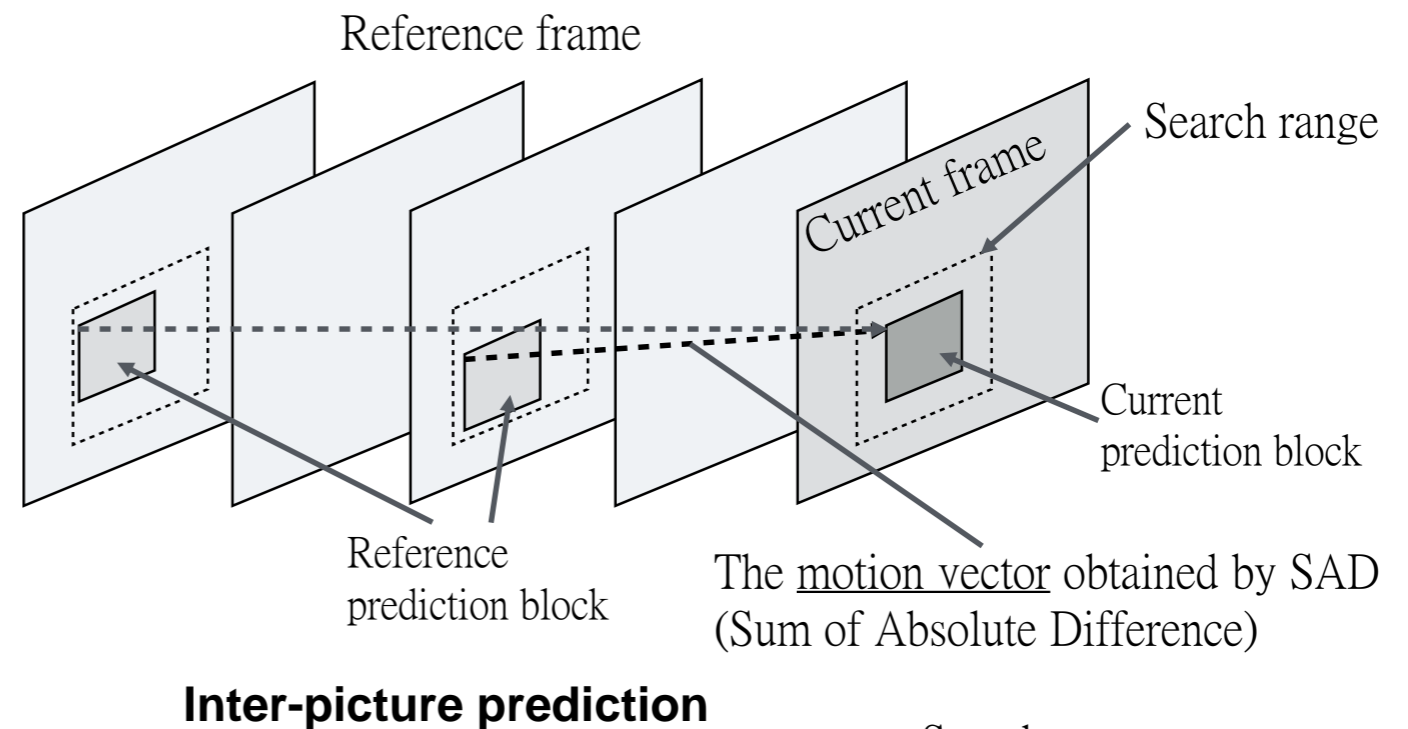

**HEVC video encoder / decoder.**

# HEVC Coding Unit Structure

**Coding Tree Unit** → **Coding Unit** → **Prediction Unit** → **Sum of Absolute Difference**

| CTU 0 | CTU 1 | CTU 2 | CTU 3 | ... | | CTU |
|-------|-------|-------|-------|-----|---|-----|
| ... | | | | | | |
| | | | | | | |
| | | | | | | CTU N-1 |

**Frame**

(64x64)

(32x32)

(16x16)

(8x8)

PU0

PU1  PU2

PU3

PU4

64
64    64
64 **Current PU**

Ref Frame 1
Ref Frame 2
Ref Frame 3
Ref Frame 4

(63+64+1)

0s

*Total SAD = #CU x #PU x  #SAD  x #RefFrame*
*        = 85   x   5   x 16384  x        4*
*        = 27,852,800 per CTU*

## Multi-Reference Frame Prediction

- Uni-predictive ME (P-frame)

- Bi-predictive ME (B-frame)



Reference frame

Search range

Current frame

Current prediction block

Reference prediction block

The motion vector obtained by SAD (Sum of Absolute Difference)

**Inter-picture prediction**

## Block Matching Motion Estimation

- Full search
  - High coding efficiency
  - High computing complexity

- Fast search
  - Low coding efficiency
  - e.g. Diamond search, TZ-search



Search range

- Introduction
- Related Work
- **System and Software**
- Proposed Methods and Implementations
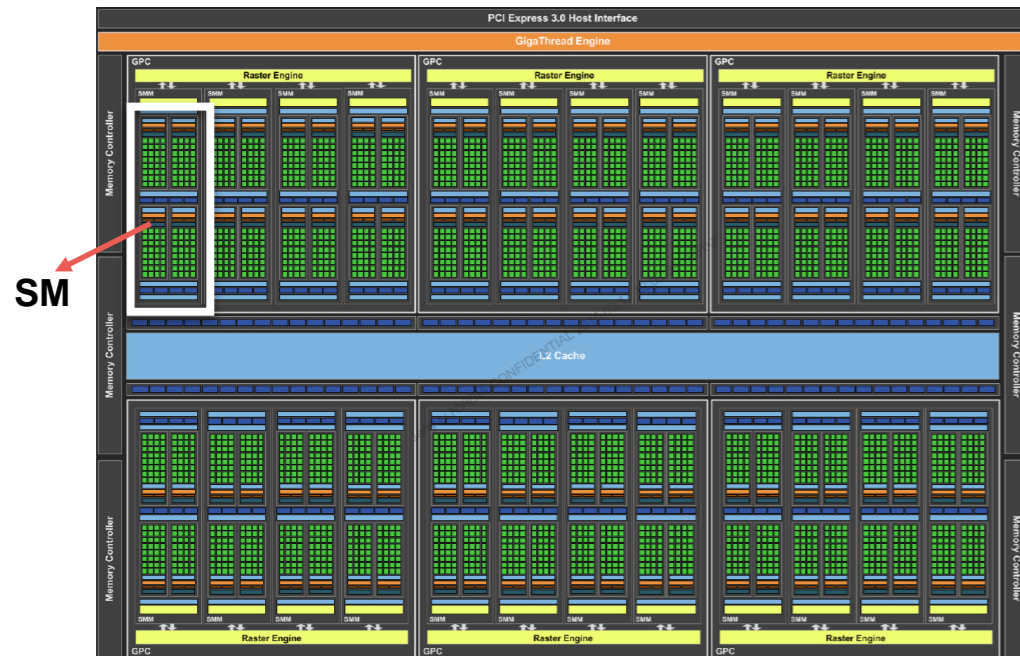- Experiment Results
- Conclusion
- Reference

# System and Software

| Machine | |
|---|---|
| **CPU** | Intel(R) Core(TM) i7-6700K, 4.00GHz, 4 cores with 8 threads |
| **GPU** | GeForce GTX **TITAN X Maxwell™** with 3072 CUDA cores, 12GB RAM |
| **Memory** | 32GB |
| **OS** | Linux ubuntu 16.04 |

| GeForce GTX TITAN X | |
|---|---|
| **Architecture** | Maxwell |
| **CUDA Capability** | 5.2 |
| **Driver Version** | 8.0 |
| **Runtime Version** | 8.0 |
| **Streaming Multiprocessors** | **24** |
| **CUDA Cores** | 3072 cores |
| **GPU Max. clock rate** | 1076 MHz |
| **Shared Memory** | 48KB/CUDA block |

| HEVC Test Model （reference software） | |
|---|---|
| **Version** | 16.9 |
| **Configure File** | encoder_lowdelay_main_P.cfg |
| **Uni-Prediction Frames** | **4 (P-Slice)** |
| **Block-Matching Algorithm** | **Full Search** |
| **Search Range** | **64** |
| **Quantization Parameter (QP)** | 22, 27, 32, 37 |
| **Frame Count** | 16 |
| **Fast Encoder Decision** | DISABLE |
| **Asymmetric Motion Partitions** | DISABLE |
| **Wavefront Parallel Processing** | DISABLE |

## GPU Hardware Architecture (Physical)



**SM**

- **24 Streaming Multiprocessors(SMs)**
- **128 CUDA cores/SM (3072 cores)**
- **12GB Global memory**

## - Streaming Multiprocessor (SM)



| Scheduler |
| 32 core | 32 core |
| 32 core | 32 core |
| Shared Mem. |

**128 cores**

**48KB programmable memory**

## CUDA Programming Model (Logical)

SIMD_CUDA_Program.cu

```
kernel<<< #CUDA block, #CUDA thread >>>(…);
{
    const Int Tid = threadIdx.x;
    const Int Bid = blockIdx.x;
    …
}
```

- **Max. number of CUDA block: $(2^{31}-1) \times (2^{16}-1) \times (2^{16}-1)$**
- **Max. number of threads per CUDA block: 1024**

| Bid | 0 | 1 | 2 | … | | | … |
|---|---|---|---|---|---|---|---|
| Kernel | 1024 thread | 1024 thread | 1024 thread | 1024 thread | … | | 1024 thread |

## - CUDA Block

- **CUDA block can run in any order**
- **Divide CUDA block into warps**

| w0 | w1 | … |
|---|---|---|
| | | |
| | | |
| | … | w31 |
| **CUDA block** | | |

**1 warp = 32 threads**

| Class | Sequence Name | Resolution | Source |
|---|---|---|---|
| **4K** | Marathon | 3840x2160 | The SJTU 4K Video Sequence Sequence Dataset **[5]** |
| | Wood | 3840x2160 | |
| | Runners | 3840x2160 | |
| | Library | 3840x2160 | |
| **A** | Traffic | 2560x1600 | HEVC Standard Standard Test Test Sequences Sequences |
| | PeopleOnStreet et | 2560x1600 | |
| **B** | BasketballDrive e | 1920x1080 | |
| | ParkScene | 1920x1080 | |
| | BQTerrace | 1920x1080 | |
| | Cactus | 1920x1080 | |
| | Kimono1 | 1920x1080 | |
| | Tennis | 1920x1080 | |

- **Introduction**
- **System and Application**
- **Proposed Methods and Implementations**
  - **Versions 1, 2, 3, 4**
- **Experiment Results**
- **Conclusion**
- **Reference**

## Sequential Execution

- Follow HM reference software



for all CU in a CTU do:
    for all PU in a CU do:
        for all RefFrame do:
            for all Pixel(x, y) in SrchRng do:

$$Distortion = SAD + \lambda_{pred} * (MVD_x + MVD_y)$$

## Proposed Methods

- Offload SAD Calculation to GPU



12

| Test Sequence | 4K: Marathon_3840x2160, 16 frames, QP=32 |
|---|---|



**Version 1**

GPU idle 74.7 sec.

**Version 2**

GPU idle 75.3 sec

**Version 3**

GPU idle 41.9 sec.

**Version 4**

GPU idle 0.79 sec.

13

## Pixel-level Parallelism

- Create 27,852,800 threads each calculating one SAD



- Memory of SAD results = 27,852,800 x sizeof(Int) = 111.4 MB per CTU
- For 4K video (2040 CTUs), it's need 227GB
- FindMv kernel to find a best MV for each prediction unit (PU)

| Test Sequence | **4K**: Marathon_3840x2160, **16** frames, QP=32 |

## CTU-level Parallelism

## CTU-level Parallelism



① **How to efficiently compute SADs using CUDA blocks?**
② **How much memory is required for each SAD kernel?**
③ **How to determine the number of threads for FindMv?**

| Test Sequence | 4K: Marathon_3840x2160, 16 frames, QP=32 |
|---|---|



5th Frame

6th Frame

CTU 0

CTU N-1

CTU 0

**CPU**

pre-pare | pre ME | ME | post ME | ... | pre ME | ME | post ME | pre-pare | pre ME | ME | post ME | ...

**GPU**

All CTUs
SAD

35.275 sec

All CTUs
SAD

time

SIMD Kernel

SAD (CTU 0 ~ CTU N-1) | F

...

SAD (CTU 0 ~ CTU N-1) | F

2.87s     1ms
C0P0 (64x64)

0.02s   1ms
C84P4 (4x8)

**Concurrent CPU & GPGPU Execution**

**① Number of CTUs per Batch?**
**② How much memory is required for each SAD kernel?**
**③ How to determine the number of threads of FindMv?**

**Multi-Threaded CPUs**

# Proposed Method Ver. 4(cont'd)

- Multi-threaded encoding based on Wave-front Parallel Processing (WPP)



HEVC WPP

**Problems:**
① **HM reference software is <u>not</u> designed for full feature multi-threading**
② **CTU-level dependency between CPU threads**
③ **How many active threads?**

| Test Sequence | 4K: Marathon_3840x2160, 16 frames, QP=32 |
|---|---|

**Version 1**



GPU idle 74.7 sec.

**Version 2**



GPU idle 75.3 sec

**Version 3**



GPU idle 41.9 sec.

**Version 4**



GPU idle 0.79 sec.

23

- Introduction
- System and Software
- Proposed Methods and Implementations
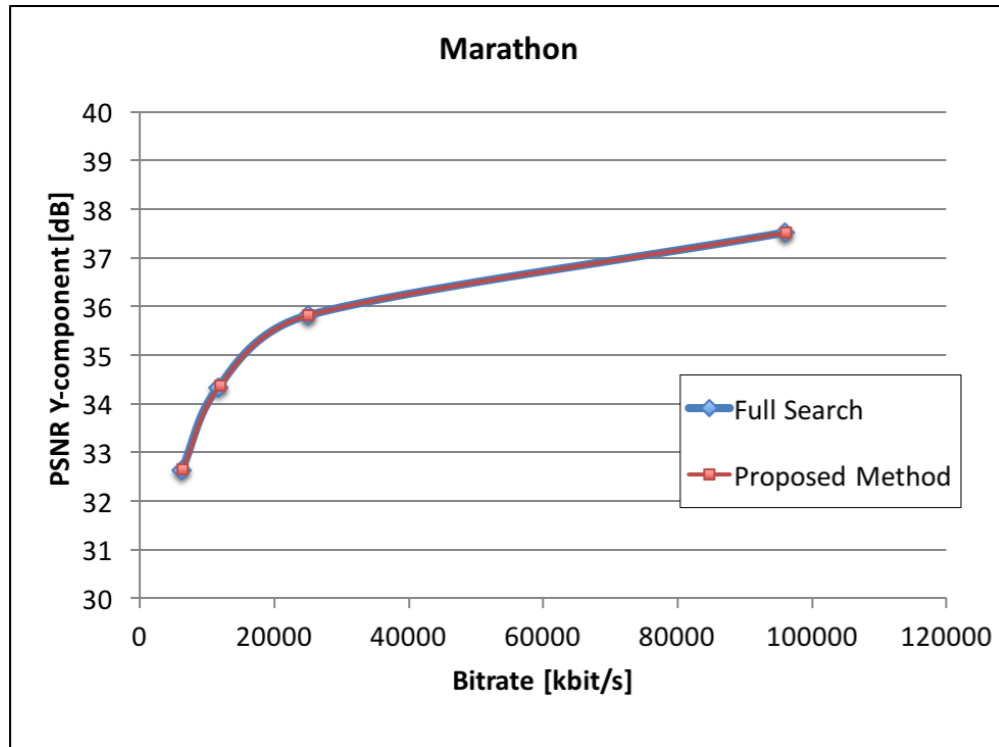- **Experiment Results**
- Conclusion
- Reference

| Version | BD-BR [%] | BD-PSNR [dB] | Motion Estimation | | Total Encoding | |
|---|---|---|---|---|---|---|
| | | | Time [s] | Speed-up | Time [s] | Speed-up |
| **Sequential** | - | - | 25666.2 | - | 26196.8 | - |
| **1** | 0.138 | -0.0037 | 550.56 | 45.41 | 1050.66 | 22.36 |
| **2** | 0.143 | -0.0039 | 255.50 | 101.50 | 756.71 | 30.74 |
| **3** | 0.143 | -0.0039 | 241.84 | 104.39 | 506.01 | 43.66 |
| **4** | 0.177 | -0.0050 | 247.79 | 102.03 | 278.94 | 90.91 |

\* *Gray column is the results of Traffic_2560x1600 sequence at QP 32.*

48.07 s

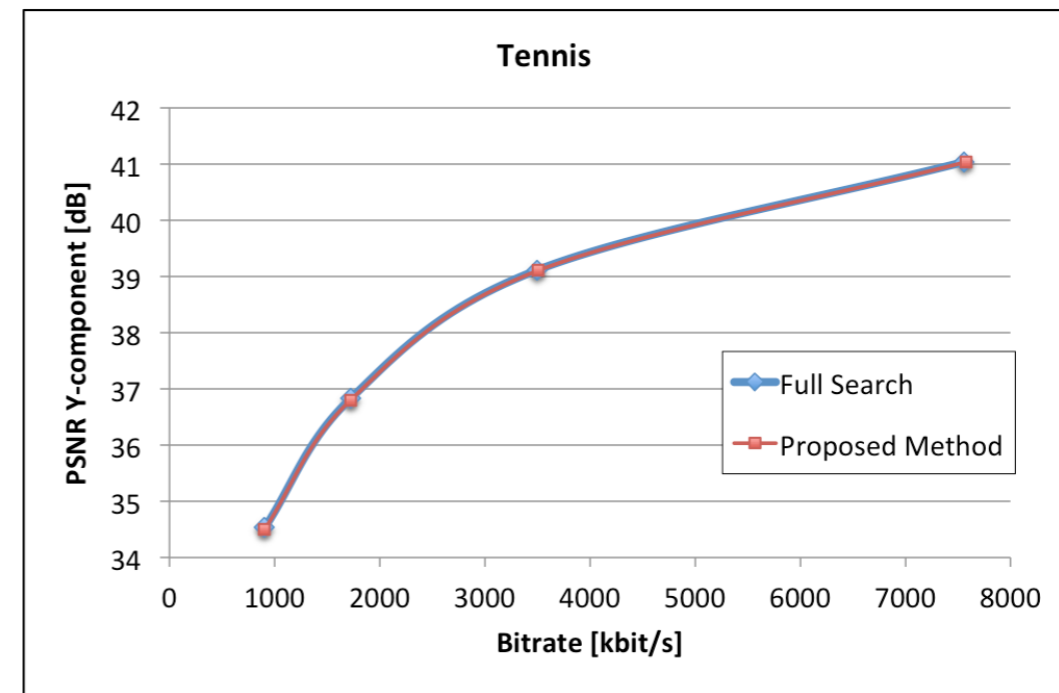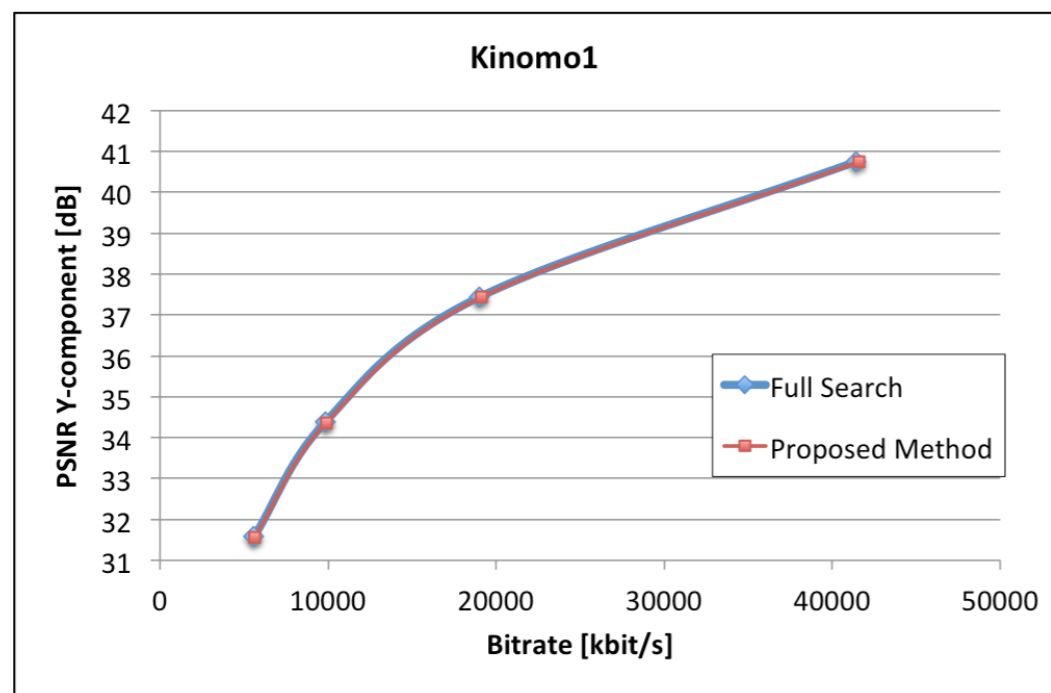The RD-curve of **4K** sequences (3840x2160)
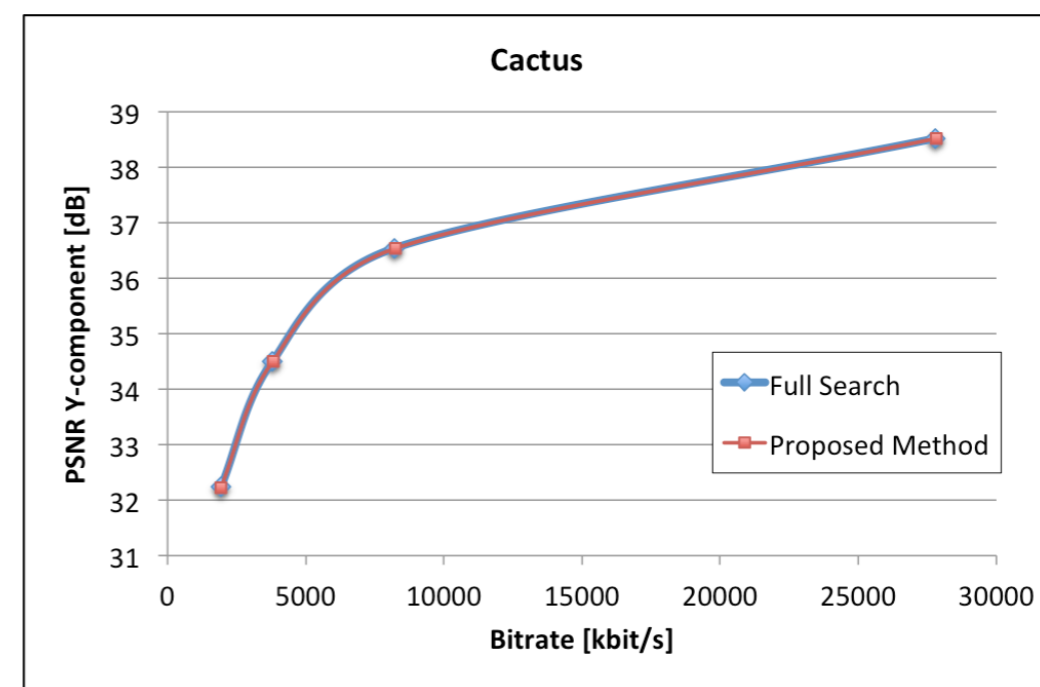
# Rate-Distortion Curves(cont'd)
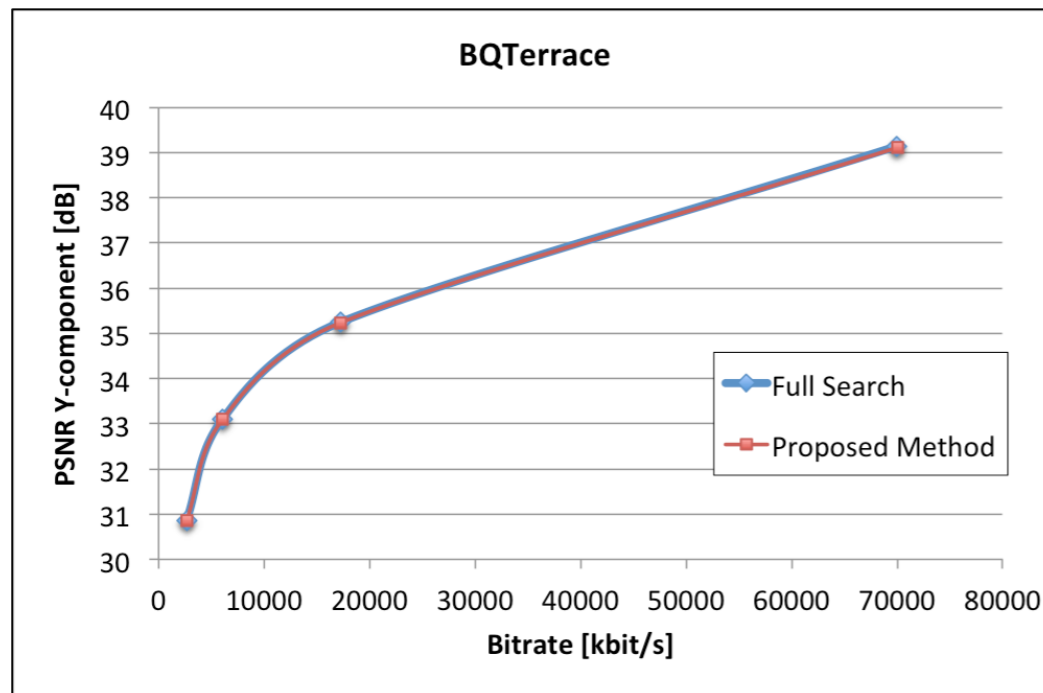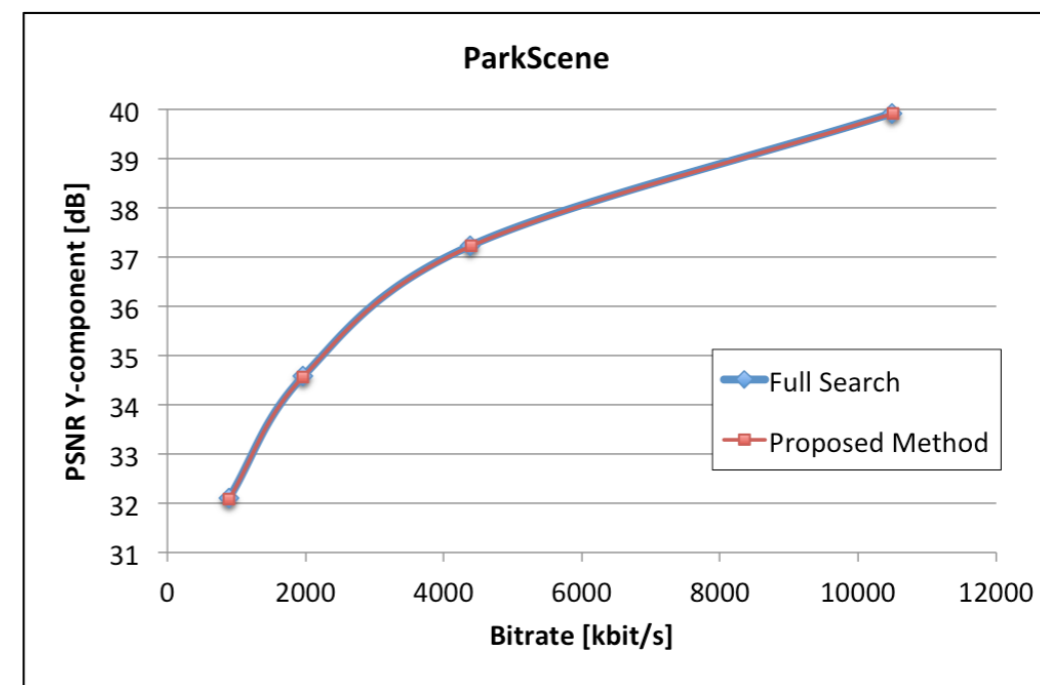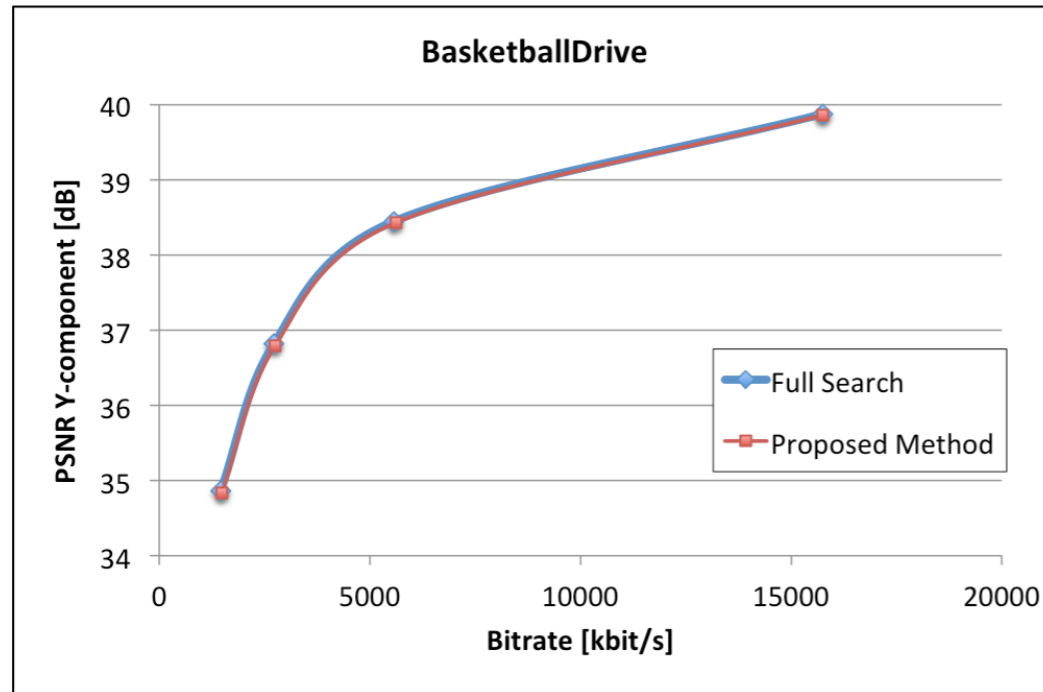


Figure. The RD-curve of class A sequences (2560x1600)



The RD-curve of **class B** sequences (1920x1080)

The RD-curve of **class B** sequences (1920x1080)

- **GPGPU and Multicore**

- **90X Speed-Up of HEVC Video Coding**

- **102X Speed-Up of Motion Estimation**

- **0.177% bit rate increase and 0.005db PSNR loss**

- **Utilization is the Key – Memory Allocation and Access**

# Thank You!!