# DSD Final Project
# Final Report

**Hardware Implementation of Pipelined RISC-V**
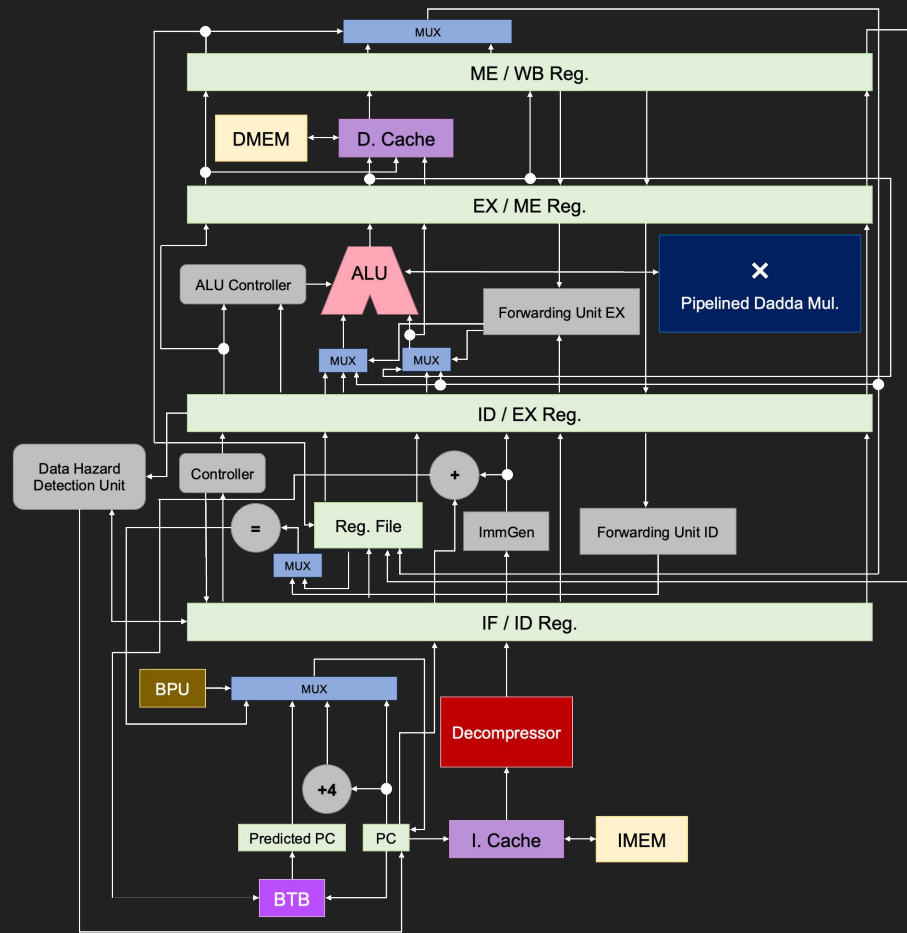
Group 7
B10901008 張禾牧
B10901016 邱巖盛
B10901078 歐信泓

# Current Results

- Complete the whole architecture.

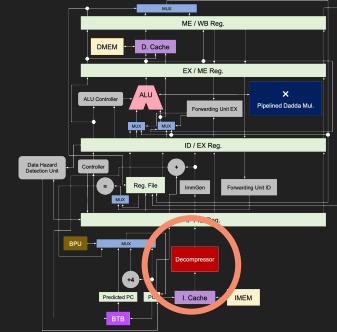- Conducted some experiments (covered later).

- **Best Performance for now**:

| $t_{cycle}$ (ns) | **2.82** |
|:---:|:---:|
| $T_{BPred}$ (ns) | 1287.330 |
| $T_{QSort}$ (ns) | 361361.85 |
| $T_{Conv}$ (ns) | 75656.37 |
| $A_{cell}$ (μm$^2$) | 313764.3873 |
| $A_{cell} T_{Conv} T_{QSort}$ | $8.58 \times 10^{15}$ |

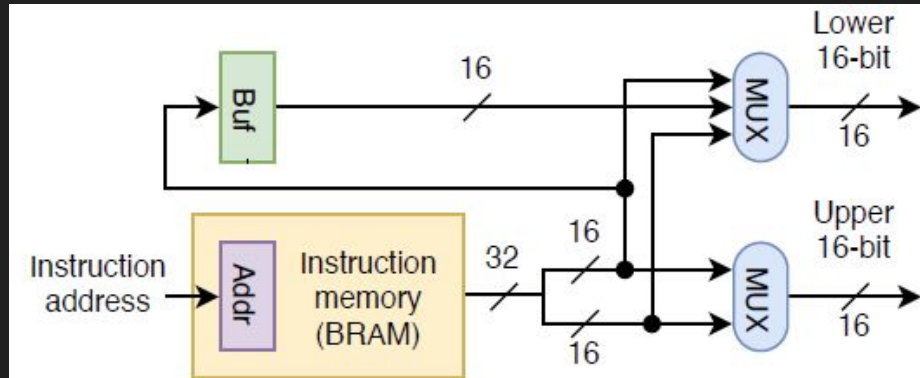# Extensions

# Decompressor

- The instruction buffer will not always be valid (e.g. after jump).

- When an invalid buffer is required (`PC[1] &&`
  `!inst_buffer_valid`), fetch lower instruction and stall for one
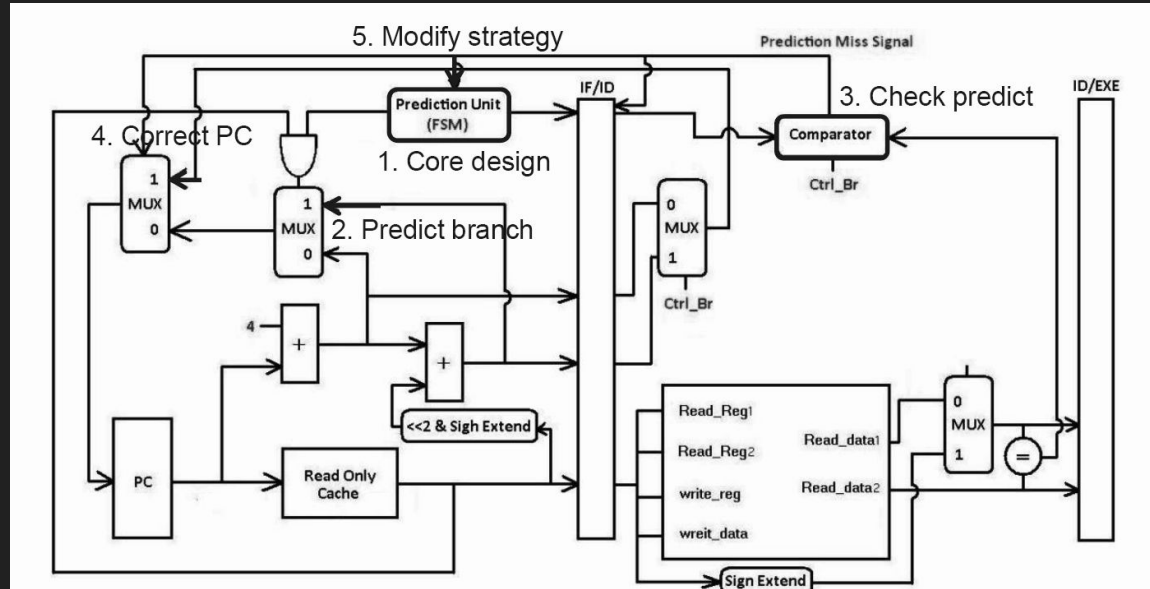  cycle

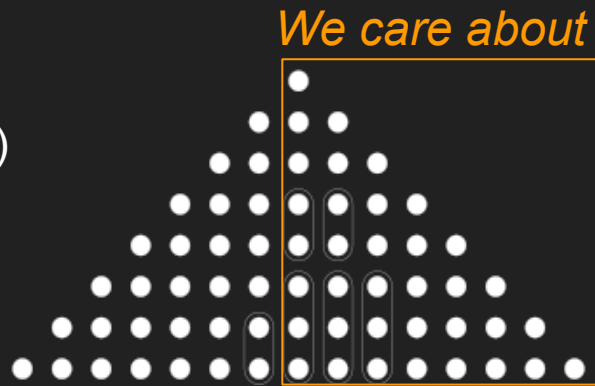- I cache addr. = `PC[31:2] + (PC[1] & inst_buffer_valid)`

# Branch Prediction

- Use two-bit saturation counter

- Save the `last_not_taken_PC` and use it if branch miss

# Multiplier

- Observations:
  - `mul` operation only requires lower 32 bits (~~mulh~~)
- Design
  - Verilog designware multiplier
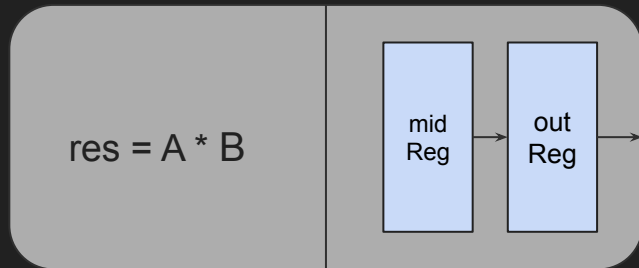  - Dadda multiplier with lower 32 bits



*We care about*

# *Pipelined* Multiplier

- Use EX & MEM (& ID) for the pipelined multiplier

- Hazard: MEM_calculation (&& ID_calculation)

- Automatic Retiming

  ```
  ○   set_dont_retime [get_registers …] false
  ○   set_optimize_register -design [get_designs mul_design]
  ```



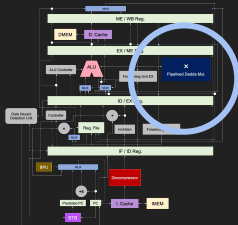res = A * B | mid Reg | out Reg

# Dadda Multiplier

First we consider the classic 8-bit multipliers:

- *dadda_8*:           basic 8-bit dadda multiplier
- *dadda_8_lower*:   last 8 bits of the product

**Partial Product Compression Priority** :
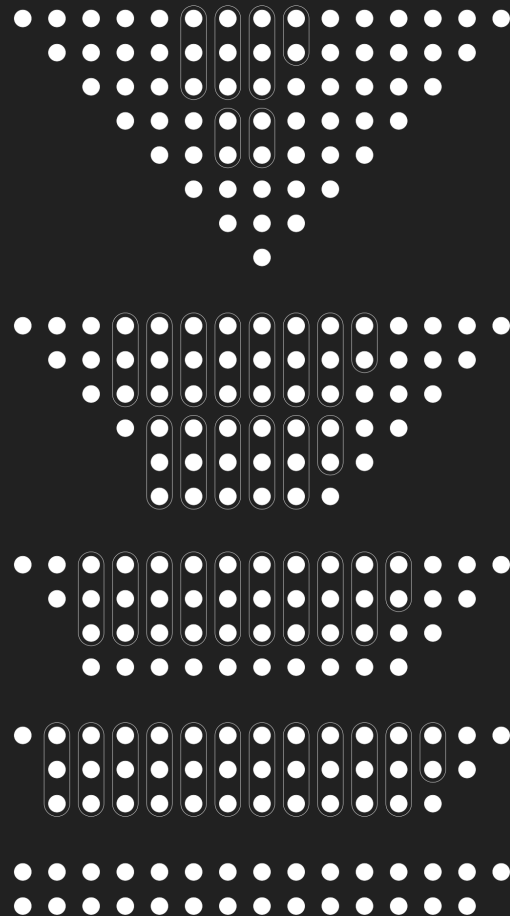
1. Data from older stage > Data from newer stage
2. Temp. carry data > Temp. sum data

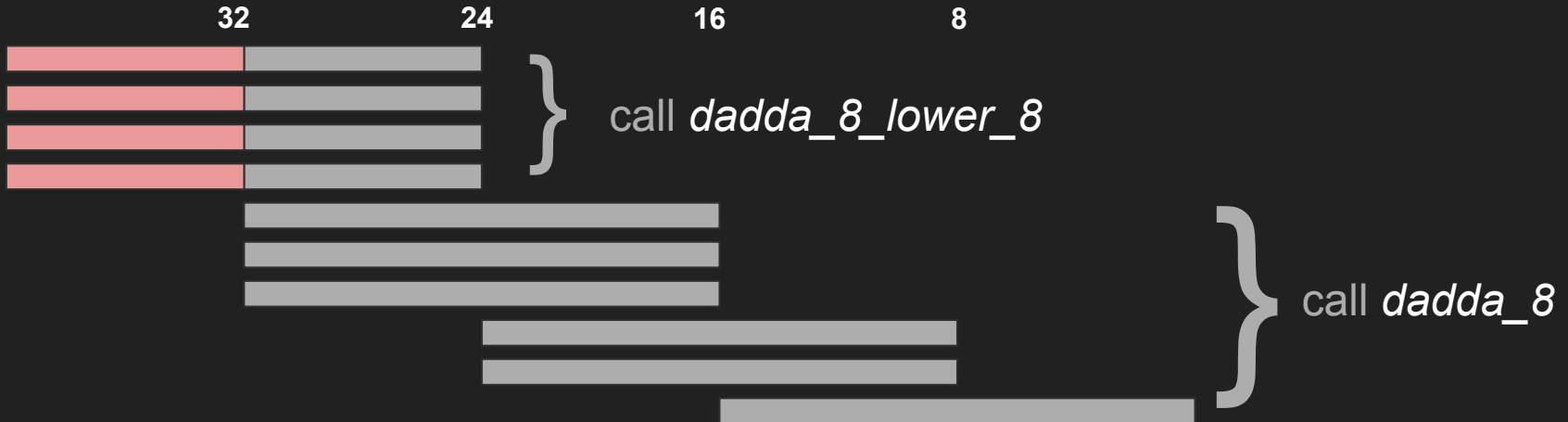    (critical path: and / or < xor )


Pipelined Dadda Mul.

# Dadda Multiplier (cont'd)

- Design a traditional 32-bit Dadda Multiplier is an extremely tough work.
- We can seperate 32-bit $a$, $b$ into concatenation of 8-bit $\{a_0, a_1, a_2, a_3\}$ and $\{b_0, b_1, b_2, b_3\}$, thus

$$ab \bmod 2^{32} = \left(a_3 2^{24} + a_2 2^{16} + a_1 2^8 + a_0\right)\left(b_3 2^{24} + b_2 2^{16} + b_1 2^8 + b_0\right) \bmod 2^{32} =$$

$$\left(a_3 b_0 + a_2 b_1 + a_1 b_2 + a_0 b_3\right) 2^{24} + \left(a_2 b_0 + a_1 b_1 + a_0 b_2\right) 2^{16} + \left(a_1 b_0 + a_0 b_1\right) 2^8 + a_0 b_0$$

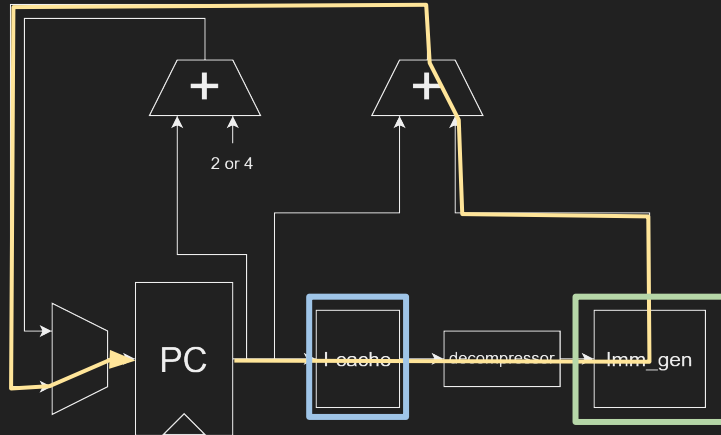call *dadda_8_lower_8*

call *dadda_8*

# Observations

# Issue about critical path

- Calculating the branch target consumes lots of time
- We use three different ways to reduce delay:
  - Instruction Cache optimization
  - Partial Immediate Generator
  - Branch Target Buffer

# About the Convolution Testbench

- Most `mul` results are not required immediately.
- More cycles can be used for `mul` calculations.

# Optimizations

# Cache Optimization

- *Direct-mapping* helps reduce critical path than *2-way* (about **0.3 ns**).

- Remove write-related function from I. Cache

  - area reduced (about 15000μm$^2$)

# Partial Immediate Generator

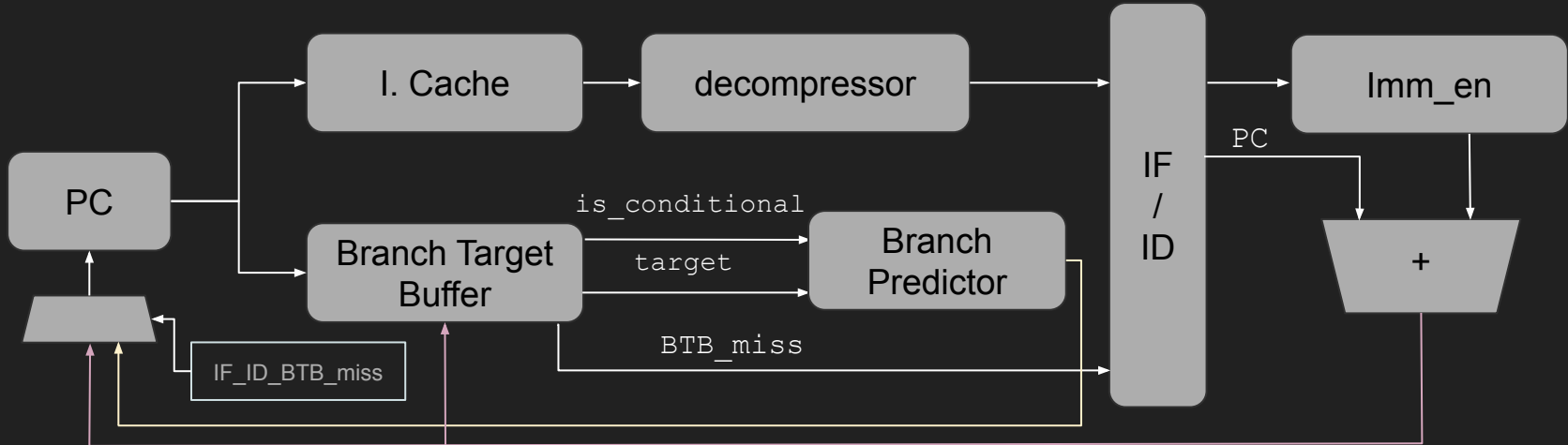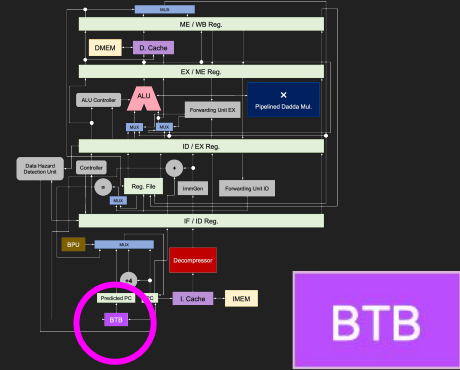- A lot of time is wasted on calculating immediate because of the *decompressor*. (decompressor → imm_gen)
- Another unit *Part_Imm_Gen* for building a faster path for branch target
  - Use "decode" technique
  - Only provide correct immediate for `(C.)beq, bne, jal`
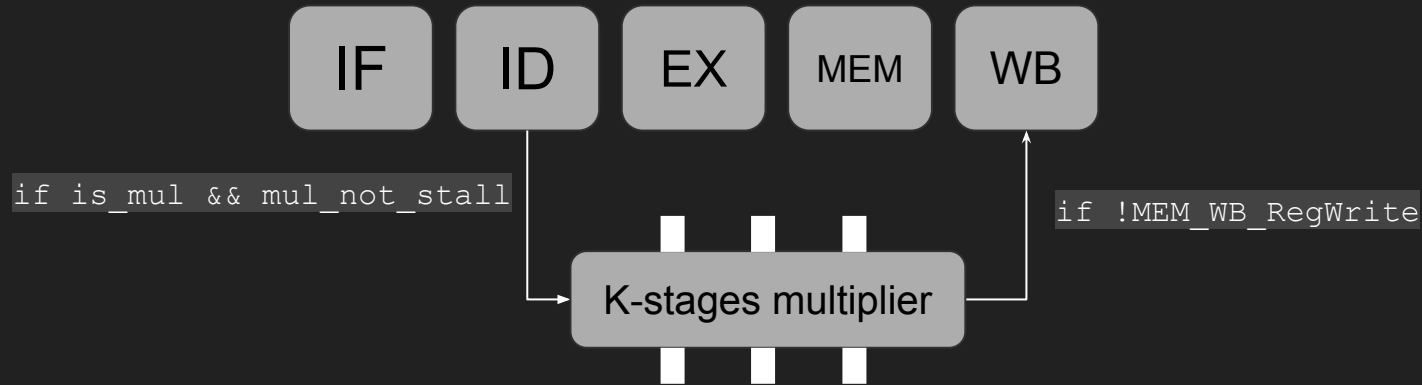- Immediate is 21 bits (partial imm is 20 bits)

# Branch Target Buffer

- Cache the target of branch instruction
- If there is any miss, the target will be calculated in ID stage & store into the buffer
- Cache Size: 2/4 block, *2-associative*

# Out of order multiplication

- It is for reducing the total area.
- No significant improvement (since the multiplier design doesn't changed).
- `mul` hazard: `[rs1 or rs2 ∈ any middle registers]`
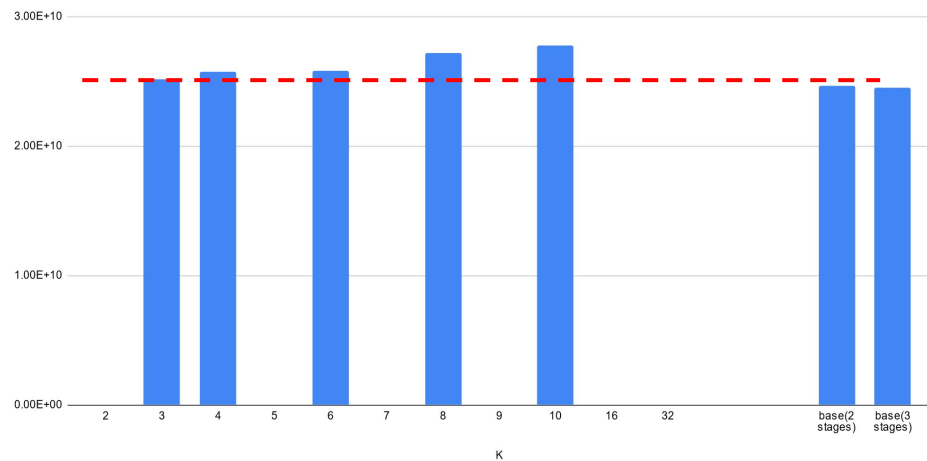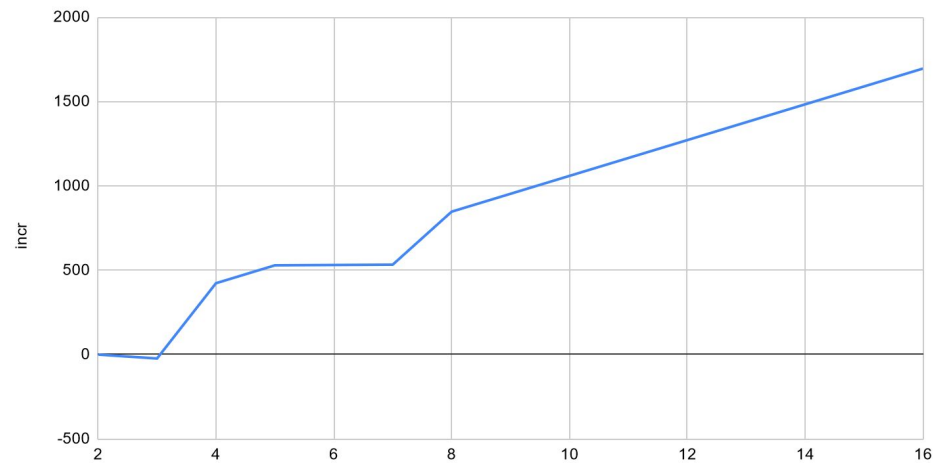- `mul_stall` = `mul_out_rd && MEM_WB_RegWrite`

# Result Analysis
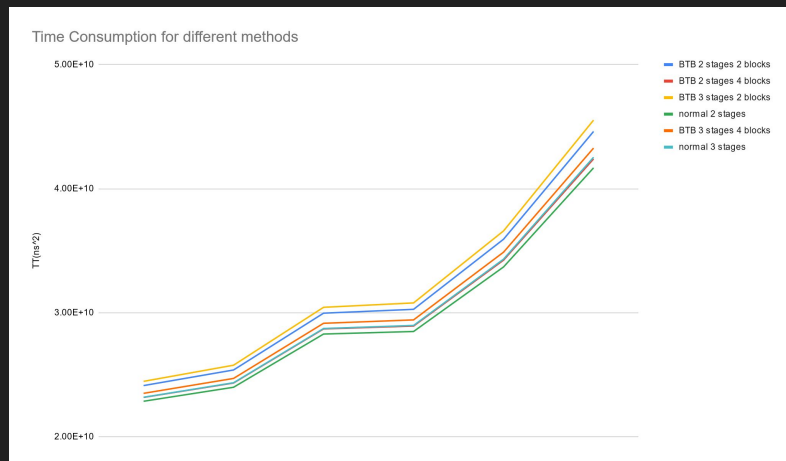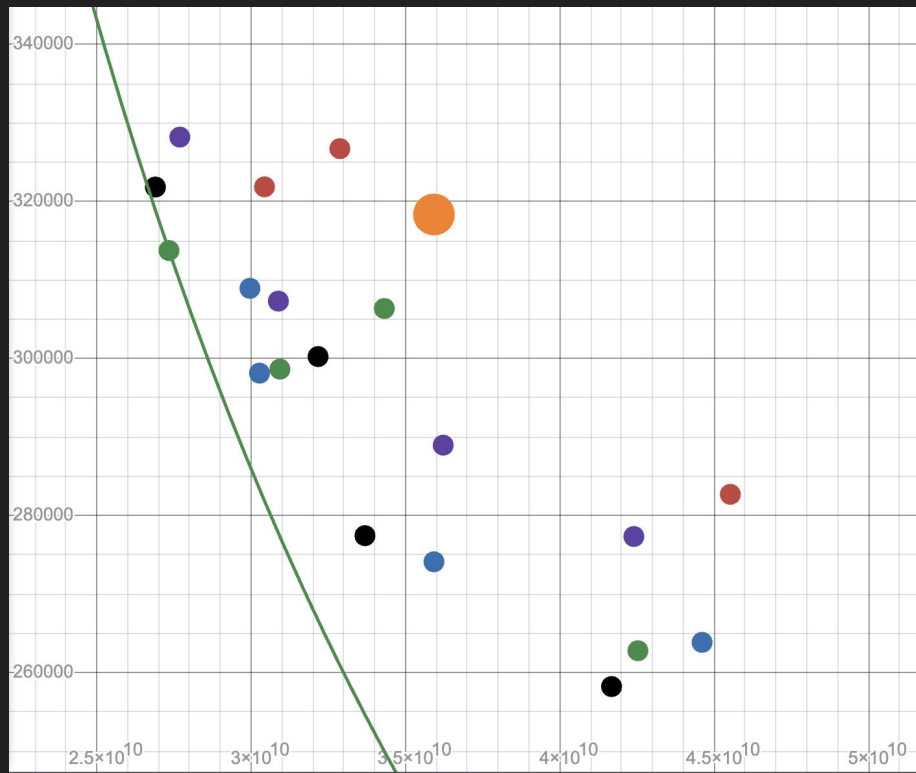
# Out-Of-Order multiplier



Area, Timing in Conv (ns) and cycles



cycle increment vs. K

# Performance



Time Consumption for different methods

- BTB 2 stages 2 blocks
- BTB 2 stages 4 blocks
- BTB 3 stages 2 blocks
- normal 2 stages
- BTB 3 stages 4 blocks
- normal 3 stages

Blue Dot:     BTB two stage(2 block)
Red Dot:      BTB three stage(2 blocks)
Black Dot:    normal two stage(2 block)
Green Dot:    normal three stage(2 blocks)
Purple Dot:   BTB two stage(4 blocks, 2 way)
Orange Dot:  Dadda multiplier

Green Curve: AT$^2$ curve passing the Green Dot

# Work Assignment Chart

- **張禾牧**：

  *Decompressor & Other Units, Optimization, Conduct Experiments*

- **邱嚴盛**：

  *Multiplier Connection, Optimization, Conduct Experiments*

- **歐信泓**：

  *Multiplier Connection, Synthesis Setup, Optimization, Conduct Experiments*