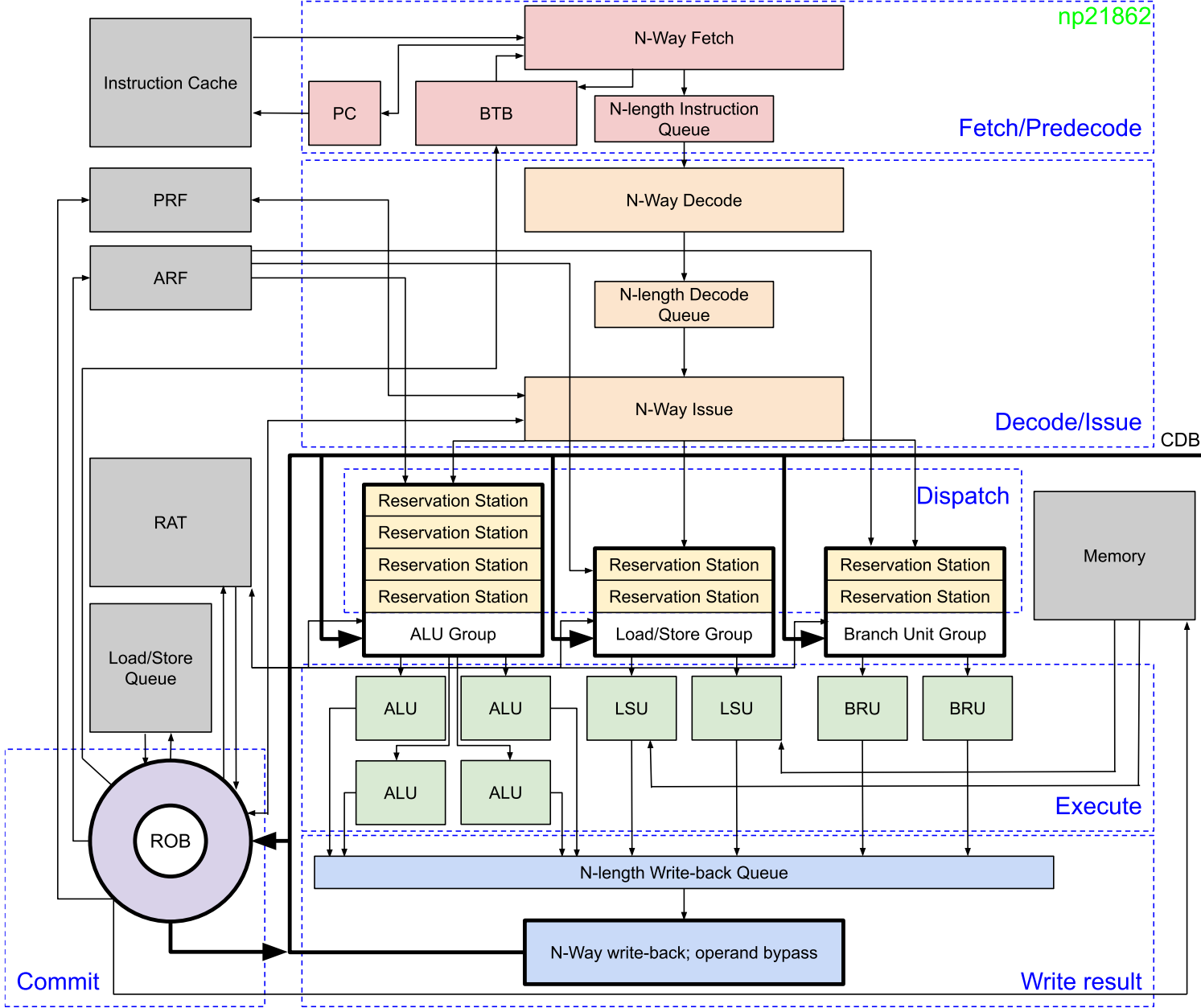


Architecture Diagram

Shown: an example configuration of my simulator with four ALU reservation stations, four ALUs, two load/store reservation stations, two load store units, two branch unit reservation stations, and two branch units.



Pipeline:

- Six-stage pipeline (fetch/predecode, decode/issue, dispatch, execute, write-result, commit)
- Out of Order execution, re-order buffer, register alias table
- Grouped reservation stations grouped by execution unit types
- Tomasulos algorithm and register renaming with the ROB
- Execution unit operand forwarding, OoO memory reads
- Aligned/unaligned instruction fetch
- Nested speculative execution with RAT history
- Dynamic branch prediction using a BTB/BTAC (1 or 2 bit predictor)
- Multi-cycle instruction latency (configurable cycle count for each instruction)
- Indexed load and store instructions
- Fixed taken & not taken and static branch backward taken & not taken predictors
- Store-to-load forwarding or store-load conflict-based flushing
- ≤ 8 ALUs, ≤ 4 LSUs, ≤ 4 branch units
- Any number (≥ 1) of reservation stations in each group

Registers and Memory (default settings):

- 256B working memory
- 32 architectural registers corresponding to those in RISC-V
- 64 physical registers
- 64 ROB entries
- 32 BTB entries
- Six different branch predictor modes

Superscalar:

- N-way fetch/decode/issue/writeback/commit per cycle
- Configurable number of ALUs, Load/Store Units, Branch Units, reservation stations

Flushing:

- All units can flush from a signaled ROB entry

Configurability & Usage:

- Command-line arguments for superscalar width, predictor type, BTB size, ALU count, load/store unit count, branch unit count, ALU reservation station count, LSU reservation station count, branch unit reservation station count, ROB size, PRF size, and an aligned-fetch flag.
- Console pipeline view (**see next slide**)

Pipeline View and Benchmark Programs

Benchmark	Information (cycle count is for single-fetch OoO)
<i>bubble_sort</i>	Sort 40 element list of integers in-place in memory using bubble sort (~17k cycles)
<i>collatz</i>	Calculate the collatz sequence from 3144 of length 61 using mod and quot routines (~98k cycles)
<i>fact</i>	Computes factorial of 14
<i>mat3_mul</i>	Stores two 3x3 matrices, then calculates and stores their product
<i>pi</i>	Computes the first few digits of pi from a grid of points and counts those < distance 1 from the origin (~300k cycles)
<i>tn</i>	'Taken not-taken' does not take branch A and takes branch B (repeats 640 times)
<i>ttn</i>	'Taken taken not-taken' takes branch A, then does not, then takes branch B (repeats)
<i>tttn</i>	'Taken x4 not-taken' takes branch A thrice, does not take it, then takes branch B (repeats)
<i>unrolled_maths</i>	A loop of mostly independent simple maths unrolled (~26k cycles)

Shown: console pipeline view of *collatz* finishing and printing register and memory content below. Run with aligned fetch, and four-way superscalar, and with a small PRF. Units are shown as a two-letter code and pipeline buffers sit between them. Instructions are shown in the pipeline by their id. Cycles is shown at the end of each like proceeding @. Reservation stations are shown as {} maps and in groups before the EU groups.

```
[FE [17,16,15,14,] DE [15,14,13,12,] IS NRS ({0=05, 1=06, 3=06, 4=05, 5=06, 6=05, 7=06}, {0=07, 1=07, 2=07, 3=07}) (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05454] pc 7
[FE [17,16,15,14,] DE [15,14,13,12,] IS NRS ({0=05, 1=06, 3=06, 4=05, 5=06, 6=05, 7=06}, {0=07, 1=07, 2=07, 3=07}) (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05455] pc 5
Flush from robEntry 96418
[FE [17,16,15,14,] DE [15,14,13,12,] IS [({,}, {,})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05456] pc 18
[FE [13,12,11,10,] DE [10,9,26,25,08,] IS [({,}, {,})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05457] pc 14
[FE [17,16,15,14,] DE [13,12,11,10,] IS [({0=25, 1=09}, {,}{0=08, 1=26})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05458] pc 14
[FE [17,16,15,14,] DE [17,16,15,14,] IS [({0=10, 1=11, 2=12, 3=13}, {,}{,})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05459] pc 14
[FE [17,16,15,14,] DE [17,16,15,14,] IS [({0=14, 1=11, 2=15, 3=13, 4=16}, {,}{0=17})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05460] pc 14
[FE [17,16,15,14,] DE [17,16,15,14,] IS [({0=14, 1=14, 2=15, 3=15, 4=16, 5=16, 6=15}, {,}{0=17, 1=17})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05461] pc 14
[FE [17,16,15,14,] DE [16,15,14,17,] IS PRF ({0=14, 1=14, 2=15, 3=15, 4=16, 5=16, 6=15}, {,}{0=17, 1=17}) (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05462] pc 14
[FE [17,16,15,14,] DE [16,15,14,17,] IS [({0=14, 1=14, 2=15, 3=15, 4=16, 5=16, 6=15}, {,}{0=17, 1=17})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05463] pc 14
[FE [17,16,15,14,] DE [16,15,14,17,] IS NRS ({0=14, 1=16, 2=14, 3=15, 4=16, 5=16, 6=15, 7=15}, {,}{0=17, 1=17, 2=17}) (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05464] pc 14
[FE [17,16,15,14,] DE [16,15,14,17,] IS NRS ({0=14, 1=16, 2=14, 3=15, 4=16, 5=16, 6=15, 7=15}, {,}{0=17, 1=17, 2=17}) (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05465] pc 14
[FE [17,16,15,14,] DE [15,14,17,16,] IS PRF ({0=16, 1=16, 2=14, 3=14, 5=16, 6=15, 7=15}, {,}{0=17, 1=17, 2=17, 3=17}) (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05466] pc 14
[FE [17,16,15,14,] DE [15,14,17,16,] IS [({0=16, 1=16, 2=14, 3=14, 5=16, 6=15, 7=15}, {,}{0=17, 1=17, 2=17, 3=17})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05467] pc 14
Flush from robEntry 96436
[FE [32,29,28,18,] DE [32,29,28,18,] IS [({,}, {,})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05468] pc 33
[FE [24,23,34,33,] DE [32,29,28,18,] IS [({,}, {,})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05469] pc 1
[FE [04,03,02,01,] DE [24,23,34,33,] IS [({0=28}, {0=32}, {0=18, 1=29})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05470] pc 5
[FE [08,07,06,05,] DE [04,03,02,01,] IS [({0=33, 1=23}, {0=32}, {0=34, 1=24})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05471] pc 25
[FE [10,09,26,25,] DE [08,07,06,05,] IS [({0=33, 1=23, 2=01, 3=02, 4=03, 5=04}, {0=32}, {0=34})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05472] pc 11
[FE [14,13,12,11,] DE [10,09,26,25,] IS [({0=05, 1=06, 4=03, 5=04}, {,}{0=34, 1=07, 2=08})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05473] pc 15
[FE [18,17,16,15,] DE [14,13,12,11,] IS [({0=05, 1=06, 2=25, 3=09, 4=03, 5=04}, {,}{0=34, 1=07, 2=26})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05474] pc 28
[FE [18,17,16,15,] DE [15,14,13,12,] IS [({0=05, 1=06, 2=25, 3=10, 5=04}, {,}{0=34, 1=07, 2=26})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05475] pc 28
Flush from robEntry 96456
[FE [32,29,28,18,] DE [32,29,28,18,] IS [({,}, {,})] (EX ___EX ___EX ___EX ___LS ___LS ___BR ___BR ___) ({,}, {,}) WB [05476] pc 35
registers (dirty): t0:1 t3:0 t4:2 t5:1 t6:1 s1:1 s2:2 a1:61 a2:1 a3:2
memory:
[00] 1572 [01] 786 [02] 393 [03] 1180 [04] 590 [05] 295 [06] 886 [07] 443 [08] 1330 [09] 665
[10] 1996 [11] 998 [12] 499 [13] 1498 [14] 749 [15] 2248 [16] 1124 [17] 562 [18] 281 [19] 844
[20] 422 [21] 211 [22] 634 [23] 317 [24] 952 [25] 476 [26] 238 [27] 119 [28] 358 [29] 179
[30] 538 [31] 269 [32] 808 [33] 404 [34] 202 [35] 101 [36] 304 [37] 152 [38] 76 [39] 38
[40] 19 [41] 58 [42] 29 [43] 88 [44] 44 [45] 22 [46] 11 [47] 34 [48] 17 [49] 52
[50] 26 [51] 13 [52] 40 [53] 20 [54] 10 [55] 5 [56] 16 [57] 8 [58] 4 [59] 2
[60] 1 [61] [62] [63]
settings: CLOCK_SPEED_MHZ=500.0
```

Superscalar width and EU Count for IPC

Hypothesis: Up to twice the Instructions Per Cycle of some benchmark programs as we go from one to eight-way superscalar. We expect a plateau as true data dependencies begin to limit instruction parallelism.

Experiment: All tests run with 64 rob entries, 8, 4, 4 reservation stations for ALUs, LSUs and branch units respectively. Performance metric of programs expressed in terms of Instructions Per Cycle, which implies programs will execute in a shorter time.

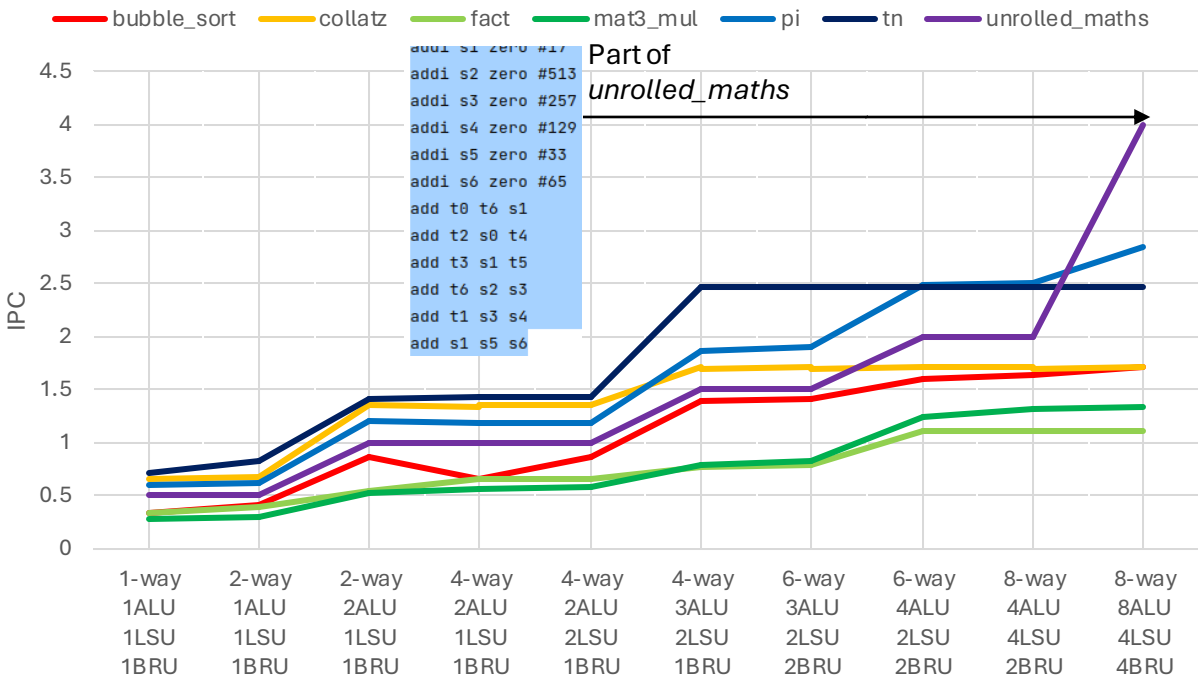
The bottom graphs shows the factor of speed-up we see from the first to the last measurement.

Result: We can see at least 4.5x IPC increase for *bubble_sort*, *mat3_mul*, *pi*, and *unrolled_maths* compared with a one-way scalar pipeline (as high as ~8x speed-up for *unrolled_maths*). The largest performance increases arise when increasing the number of ALUs in the processor, from 1 to 2 ALUs, from 2 to 4 and 4 to 8. After 4 ALUs and 6-way superscalar, for all but two programs' IPC plateaus likely due to limited instruction-level parallelism, implying data dependencies.

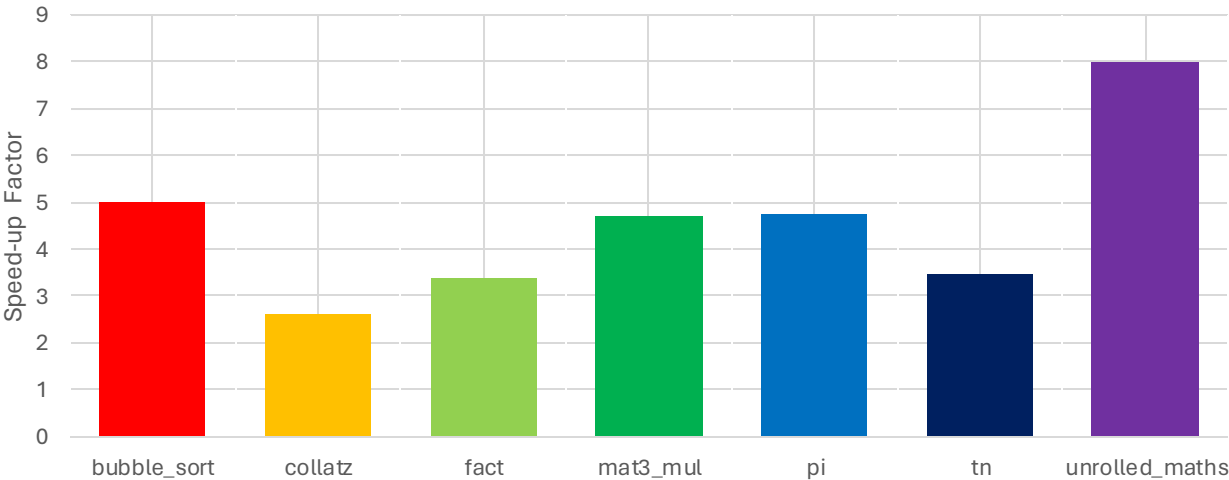
unrolled_maths and *pi* increased in IPC approx. proportionally to the number of ALUs up to 4. They contain almost no memory operations (and *unrolled_maths* contains no branches). *pi* cant keep up with *unrolled_maths* for 8 ALUs due to true data dependencies.

Bubble sort exhibits performance dips when it is allowed to speculate incorrectly for longer. Overall we accept the hypothesis.

IPC versus CPU Configuration



Ratio of Last over First IPC for CPU Configurations



Reservation Stations versus Instruction Misprediction

Hypotheses: 1) There will be at least three times the mispredicted instruction percentage for our benchmarks between lowest and higher # reservation stations, 2) programs will all see increased IPC regardless.

Experiment: Percentage of mispredicted instructions is calculated as:

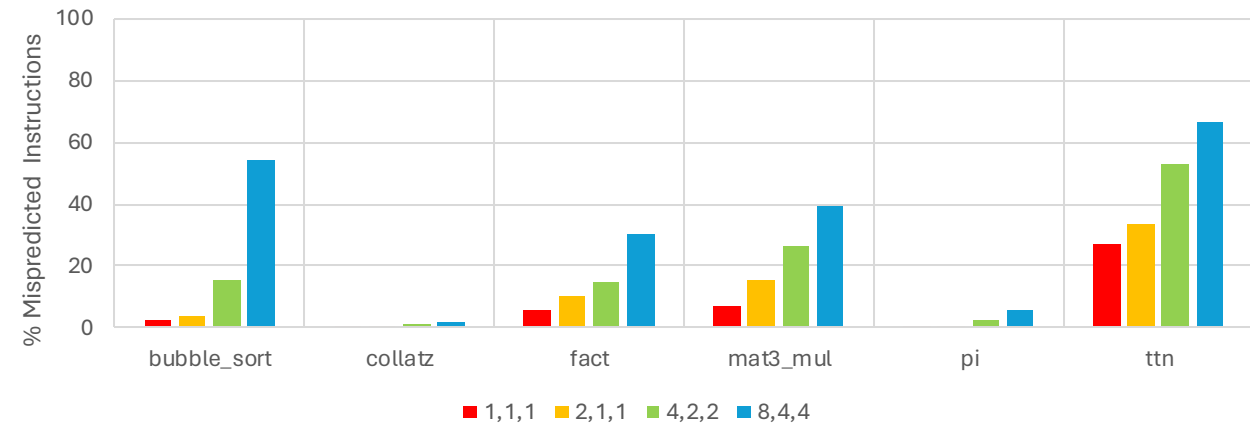
$$\frac{\# \text{ Instructions Flushed from the ROB}}{\# \text{ Committed Instructions} + \# \text{ Instructions Flushed from the ROB}}$$

Test run with 64 rob entries, 4 ALUs, 4 LSUs, 2 branch units, 8-way superscalar, two-bit predictor, four times PRF as ARF. Graphs display 3 numbers corresponding to the number of reservation stations in each group (**ALU, LSU, BRU**).

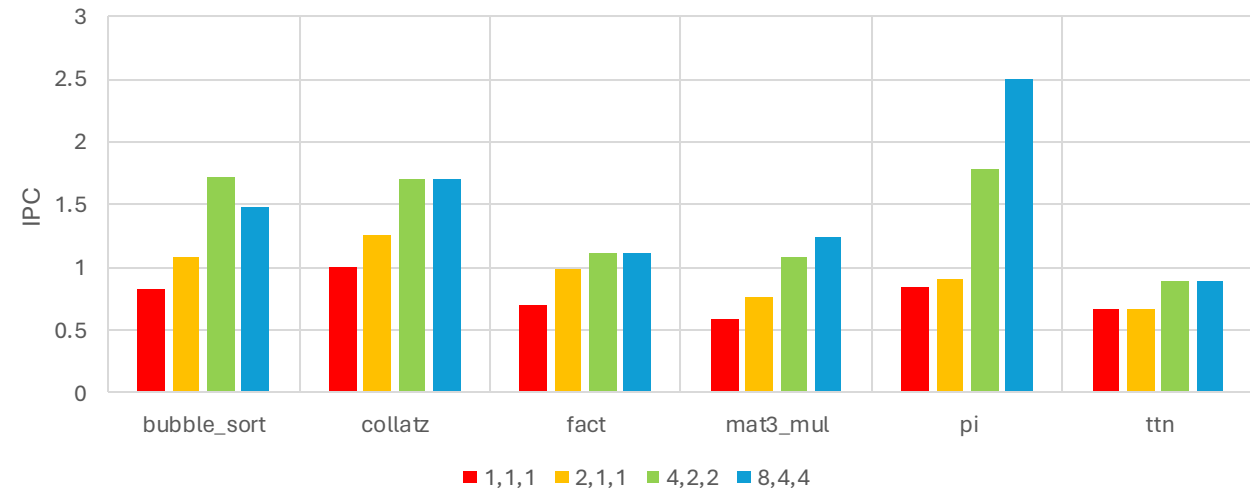
Result: We see a monotonic increase in program IPC for all programs when we increase the number of reservation stations, except for *bubble_sort* and *collatz*. Confirming our hypothesis, we see more than three times the mispredicted instruction percentage for all programs except *ttn*. Overall we accept the hypotheses, noting the exceptions.

bubble_sort loses IPC when going to 8 ALU-reservation-stations due to the wasted time and resources processing the increased number of incorrectly speculated instructions. Note percent of mispredicted instructions seems to scale approximately linearly with our exponential increase in reservation stations, indicating they allow for more instruction-level parallelism.

Increase in Percent Mispredicted Instructions against Res. Stat.ns



Variation in IPC against Num. of Reservation Stations



Hypothesis: using more sophisticated branch predictors will show an increased IPC indicative of accurate speculation.

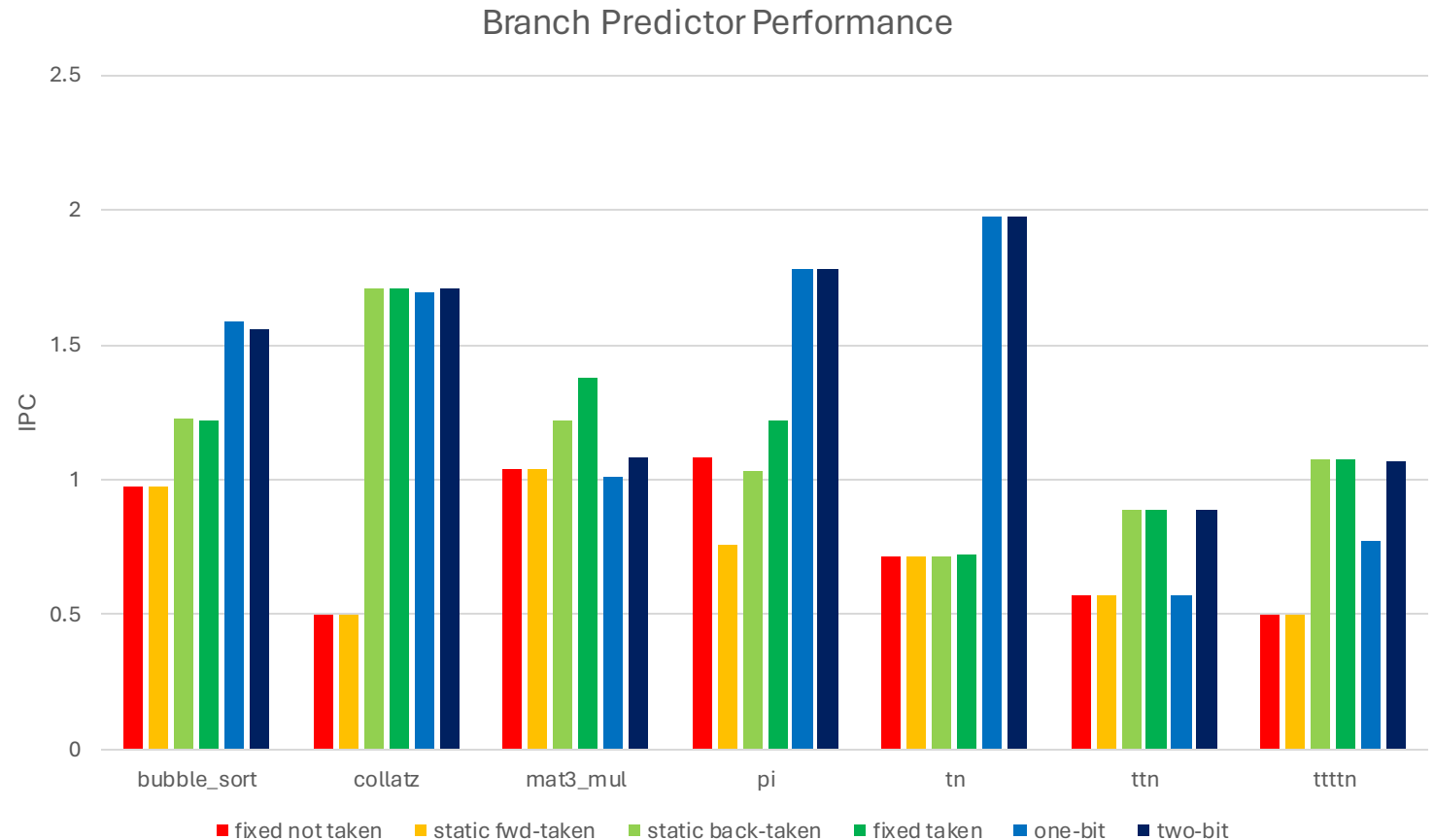
Experiment: Run 8-way superscalar with 4 ALUs, 4 LSUs, 2 Branch Units. 4,2,2 reservation stations for ALUs, LSUs and branch units resp. Showing the effect of fixed, static and one/two-predictors. The latter two use a BTB/BTAC to which also updates a state-machine for each recorded branch.

Result: Static forward-taken backward-not-taken performed the worst of all. This makes sense because I write assembly in the opposite way.

Fixed not-taken does poorly, but better than the previous in the case of *pi* because it will at least correctly predict forward branches generally aren't taken.

One and two-bit generally show the best IPC. Two-bit pulls ahead of one-bit for *ttn* and *tttn* because it quickly produces a strongly-taken state in the retaken branch's BTB entry.

In all but one case the one- or the two-bit do no worse than the others on which we accept the hypothesis.



The effect of PRF Size on IPC

Hypothesis: Limiting the size of the physical register file will cause the issue stage to block often as renaming cannot occur until physical registers are available. This will create stalls and we'll see reduced IPC by at least two-times.

Experiment: The simulator was run 8-way superscalar and sufficient headroom for high ILP. PR count is increased exponentially to investigate whether any programs require a large PRF.

Result: All programs show reduced IPC at two and four physical registers. For two, only one instruction is allowed to flow through at a time as there must be at least two free physical registers to place a register into a reservation station.

All programs drastically increase in performance from 2 to 8 physical registers. Most increase and plateau at 16. Notable outliers include *bubble_sort* and *mat3_mul* where less blocking causes more incorrectly speculated instructions and lower IPC.

pi is able to make use of the extra 16 physical registers gained in the jump from 16 to 32. To get maximum ILP, in-flight instructions have more diverse destination registers than other benchmarks. The graph clearly supports our hypothesis.

