# CPU microarchitectures (2000-2018)

# CPU Microarchitecture

- Branch prediction
  - conditions, indirect branches, call/return pairs
  - speculative execution

- Instruction decoder
  - loop cache (simple loops up to ~20 instructions)
  - conversion to micro-ops (1:1, 1:N, 2:1)
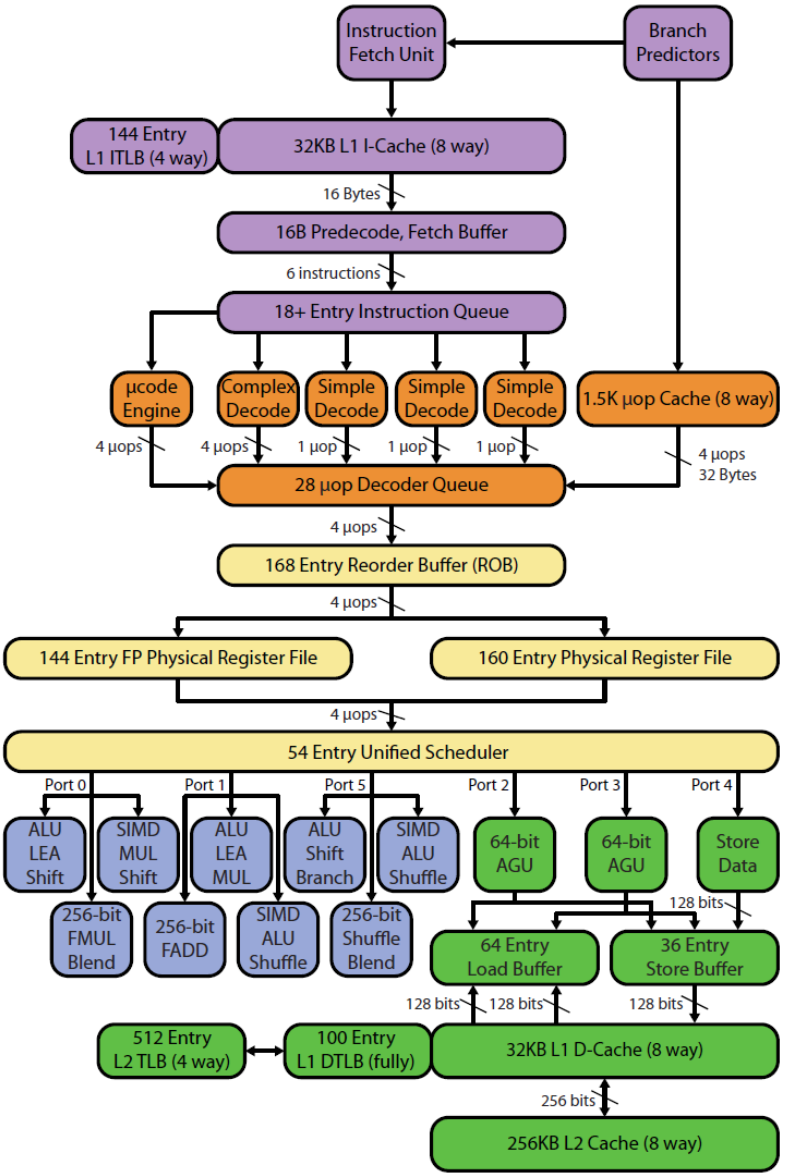  - stack-pointer simulator

- Renamer
  - 16 architectural integer registers mapped to ~160 physical
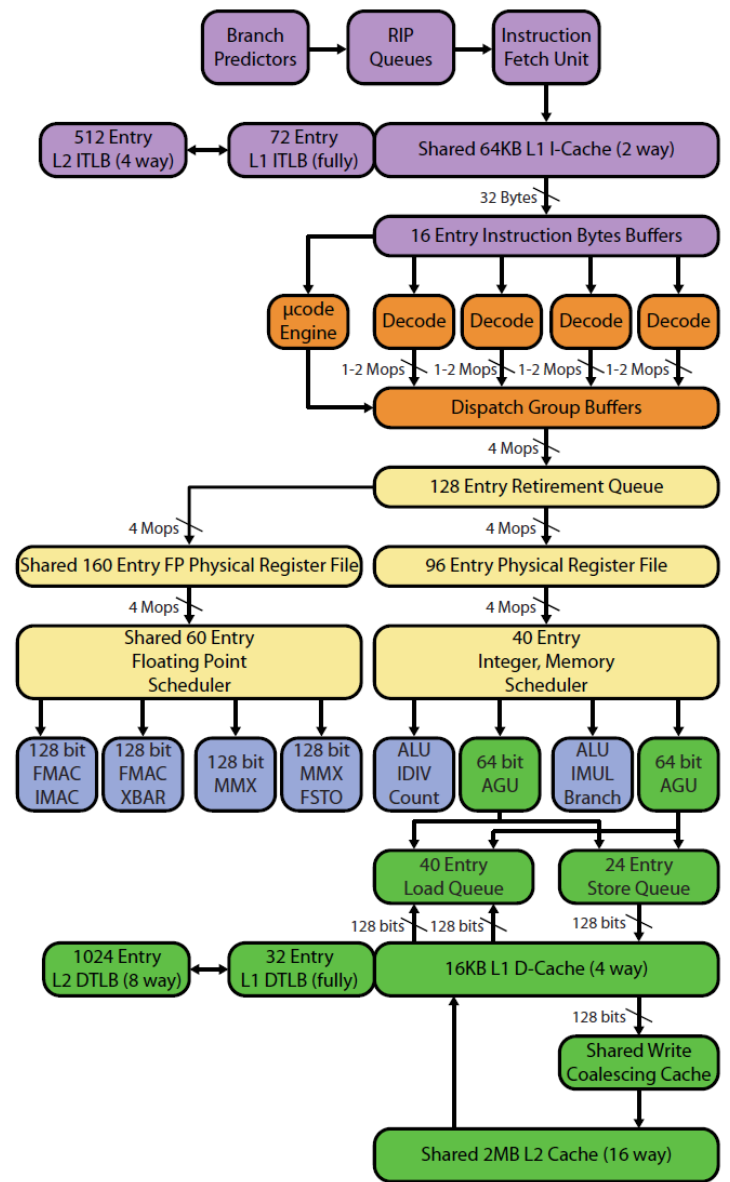    - similarly for FP/vector registers

- Out-of-order execution
  - ~40 micro-ops in simultaneous execution (RS) from a window of ~300 (ROB)
  - retirement: memory/register stores in-order in background
  - store forwarding: loads retrieve values from waiting stores
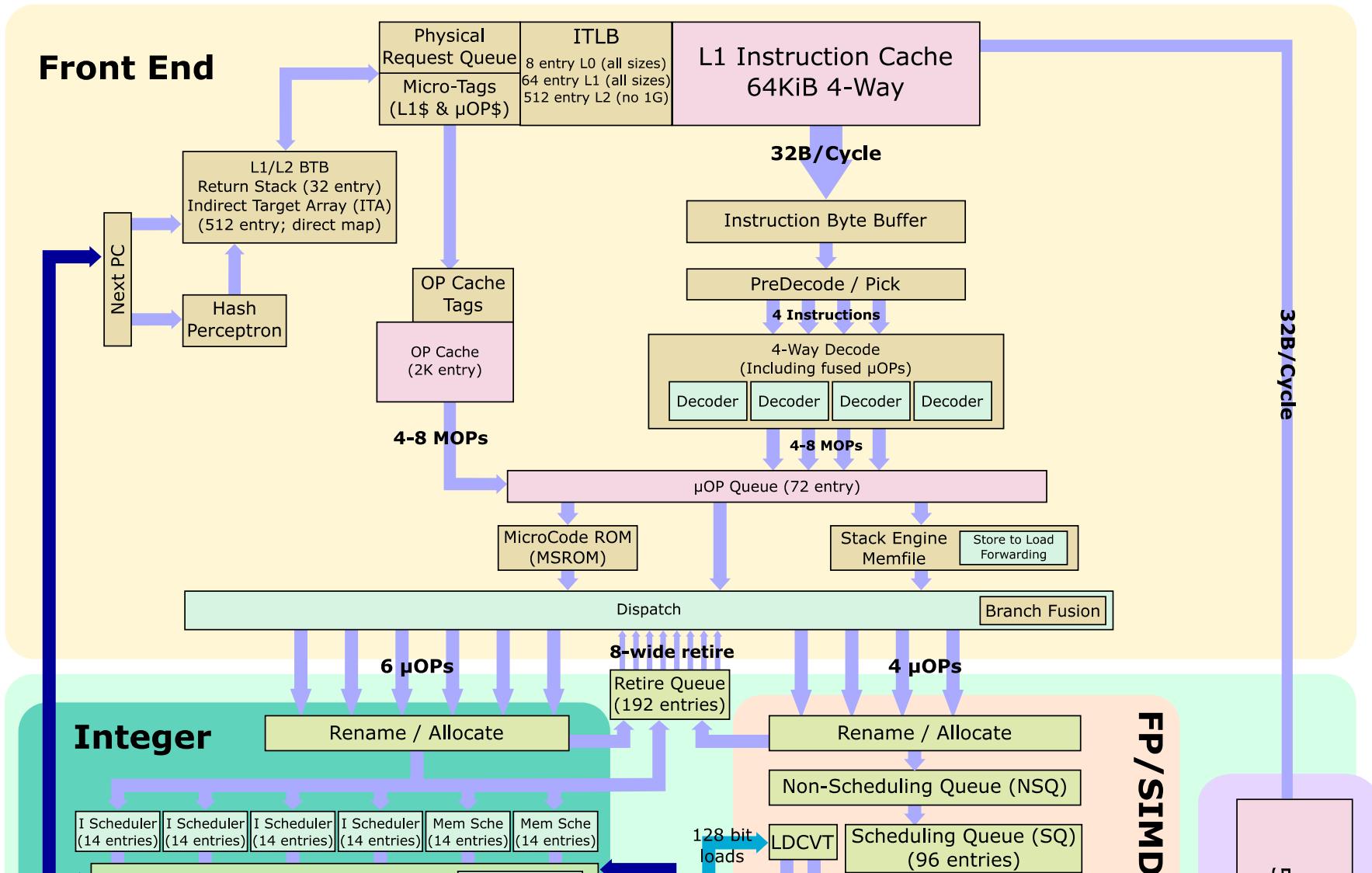  - speculative loads: no waiting for waiting stores (to unknown addresses)

# CPU Microarchitecture (Intel x86/64, 2015..2023)
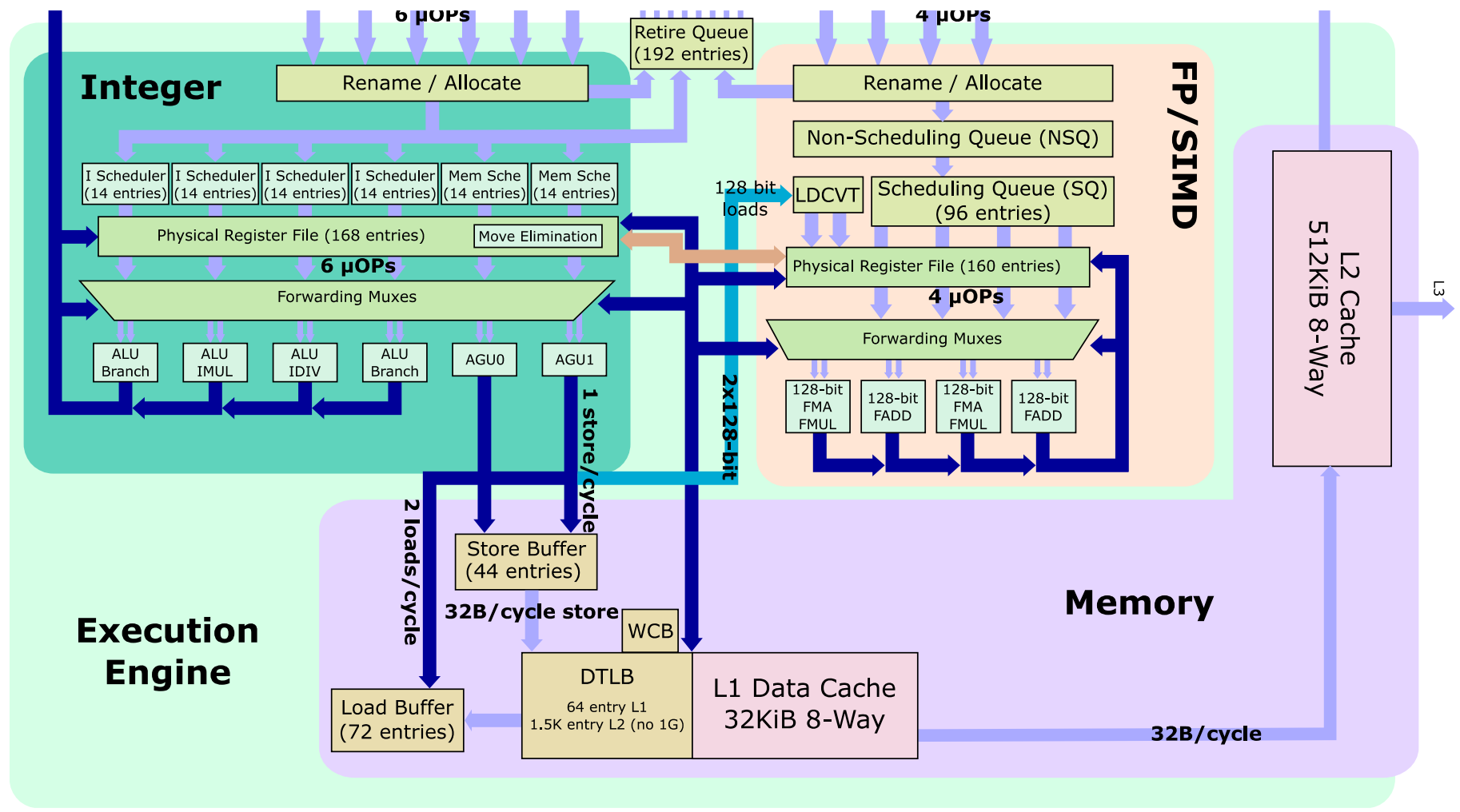
▸ In a cycle, CPU can (in theory) simultaneously perform:

  ▸ Fetch: 16 B (avg. 4 instructions) from L1 instruction cache

  ▸ Decode: 5..6 instructions

  ▸ ALU: 3..5 simple operations (add/mul)

  ▸ Memory load: 2..3 reads (up to 256..512 bits total) from L1 data cache

  ▸ Memory store: 1..2 writes (up to 256..512 bits total) to L1 data cache

▸ Latency

  ▪ the time between consuming operands and producing results

  ▸ integer add: 1, mul: 4

  ▸ FP add: 2..4, FP mul/FMA: 4

  ▸ FP div: 11/14, integer div: ~18, data dependent

  ▸ integer load: 5, FP load: 6 (L1 cache)

  ▸ store address: 3

  ▸ store data:  1..2 (retirement, in-order)