# Modern Processors

I asked you to watch the video by Cliff Click on modern hardware:

http://www.infoq.com/presentations/click-crash-course-modern-hardware

Cliff Click said that 5% miss rates dominate performance. Let's look at why. I looked up a characterization of the SPEC CPU2000 and CPU2006 benchmarks[1].

Here are the reported cache miss rates[2] for SPEC CPU2006.

| | |
| --- | --- |
| L1D | 40‰ |
| L2 | 4 ‰ |

Let's assume that the L1D cache miss penalty is 5 cycles and the L2 miss penalty is 300 cycles, as in the video. Then, for every instruction, you would expect a running time of, on average:

$$1 + 0.04 \times 5 + 0.004 \times 300 = 2.4.$$

Misses are expensive!

**Forcing branch mispredicts.** It takes a bit of trickery to force branch mispredicts. gcc extensions allow hinting, but usually gcc or the processor is smart enough to ignore bad hints. This[3] worked last year, though:

```
#include <stdlib.h>
#include <stdio.h>

static __attribute__ ((noinline)) int f(int a) { return a; }

#define BSIZE 1000000
int main(int argc, char* argv[])
{
  int *p = calloc(BSIZE, sizeof(int));
  int j, k, m1 = 0, m2 = 0;
  for (j = 0; j < 1000; j++) {
    for (k = 0; k < BSIZE; k++) {
      if (__builtin_expect(p[k], EXPECT_RESULT)) {
        m1 = f(++m1);
      } else {
        m2 = f(++m2);
      }
    }
  }

  printf("%d, %d\n", m1, m2);
}
```

---

[1] A. Kejariwal et al. "Comparative architectural characterization of SPEC CPU2000 and CPU2006 benchmarks on the Intel Core 2 Duo processor", SAMOS 2008.

[2] ‰ is "permil", or per-1000.

[3] Source: blog.man7.org/2012/10/how-much-do-builtinexpect-likely-and.html.

Running it yielded:

```
plam@plym:~/459$ gcc -O2 likely-simplified.c -DEXPECT_RESULT=0 -o likely-simplified
plam@plym:~/459$ time ./likely-simplified
0, 1000000000

real 0m2.521s
user 0m2.496s
sys 0m0.000s
plam@plym:~/459$ gcc -O2 likely-simplified.c -DEXPECT_RESULT=1 -o likely-simplified
plam@plym:~/459$ time ./likely-simplified
0, 1000000000

real 0m3.938s
user 0m3.868s
sys 0m0.000s
```

`gcc` seems to have gotten smart enough to reject bogus hints in the interim.

# Limits to parallelization

I mentioned briefly in Lecture 1 that programs often have a sequential part and a parallel part. We'll quantify this observation today and discuss its consequences.

**Amdahl's Law.** One classic model of parallel execution is Amdahl's Law. In 1967, Gene Amdahl argued that improvements in processor design for single processors would be more effective than designing multi-processor systems. Here's the argument. Let's say that you are trying to run a task which has a serial part, taking fraction $S$, and a parallelizable part, taking fraction $P = 1 - S$. Define $T_s$ to be the total amount of time needed on a single-processor system. Now, moving to a parallel system with $N$ processors, the parallel time $T_p$ is instead:

$$T_p = T_s \cdot (S + \frac{P}{N}).$$

**As $N$ increases, $T_p$ is dominated by $S$, limiting potential speedup.**

We can restate this law in terms of speedup, which is the original time $T_s$ divided by the sped-up time $T_p$:

$$\text{speedup} = \frac{T_s}{T_p} = \frac{1}{S + P/N}.$$

Replacing $S$ with $(1 - P)$, we get:
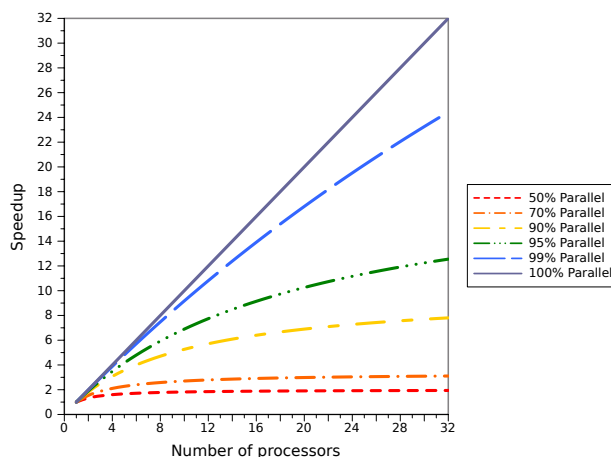
$$\text{speedup} = \frac{1}{(1 - P) + P/N},$$

and

$$\text{max speedup} = \frac{1}{(1 - P)},$$

since $\frac{P}{N} \to 0$.

**Plugging in numbers.**   If $P = 1$, then we can indeed get good scaling; running on an $N$-processor machine will give you a speedup of $N$. Unfortunately, usually $P < 1$. Let's see what happens.

| $P$ | speedup $(N = 18)$ |
|---|---|
| 1 | 18 |
| 0.99 | $\sim 15$ |
| 0.95 | $\sim 10$ |
| 0.5 | $\sim 2$ |

Graphically, we have something like this:



Amdahl's Law tells you how many cores you can hope to leverage in an execution given a fixed problem size, if you can estimate $P$.

**Consequences of Amdahl's Law.**   For over 30 years, most performance gains did indeed come from increasing single-processor performance. The main reason that we're here today is that, as we saw in the video, single-processor performance gains have hit the wall.

By the way, note that we didn't talk about the cost of synchronization between threads here. That can drag the performance down even more.

**Amdahl's Assumptions.**   Despite Amdahl's pessimism, we still all have multicore computers today. Why is that? Amdahl's Law assumes that:

- problem size is fixed (read on);

- the program, or the underlying implementation, behaves the same on 1 processor as on $N$ processors; and

- that we can accurately measure runtimes—i.e. that overheads don't matter.