

INFOMOV 2019 FINAL EXAM – EDUC-THEATRON – 17:00 – 19:00

1. The 'External Memory Model' assumes that the cost of running an algorithm is proportional to the number of memory blocks that need to be read into the cache during the execution of the algorithm. Describe one scenario in which this assumption is valid, and one scenario in which this assumption is invalid, or at least mostly invalid. Valid scenarios are plenty. Invalid scenarios: anything that loads some data and operates on it endlessly. The Fibonacci example that several brought up is excellent. Not acceptable: just making things more expensive (sin/cos/sqrt on a stream of data). In that case things are still proportional to the number of accessed blocks.

2. Consider the following data structure:

```
struct Nuke
{
    float x, y, z;           // position
    bool homing;            // flag
    float vx, vy, vz;       // missile velocity
    bool exploded;          // flag
};
```

Basic idea: you want to minimize memory cost, i.e. the # of touched cache lines. For random access this means aligning the struct to cache line boundaries (make it 2^N in size). For sequential access this means making it as small as possible. Layouts for SIMD processing also yielded full points for sequential access.

Rewrite this structure twice: once for efficient random access of a large array of Nukes, and once for efficient sequential access of a large number of Nukes in a single, continuous array. Motivate your layouts. Some pitfalls for Q2: a float is automatically aligned to 4 bytes. A cacheline is not 128bit. A bool is not 1 bit.

3. Some CPU architecture questions, based on "Modern Microprocessors - a 90 minute guide":

In 30 words or less,

- a) explain what a *pipeline latch* is. Skipped by most of you, but easily extracted from the document. Quite interesting actually.
- b) explain what a *bypass* is. Still a recommended read. 😊
- c) explain what *speculative execution* is.
- d) explain what *predication* is.

4. The Core i7-8700K processor uses 2x6MB of 16-way set-associative L3 cache. AMD's K10 uses 6MB of 48-way set-associative cache.

- a) Under what circumstances is 48-way better than 16-way, and when and/or how is 16-way better?
- b) Both processors use much lower set associativity for L1 and L2. Why do you think this is?

The whole class is cache expert, but for completeness:

48-way is closer to fully-associative and reduces collisions better than 16-way. It is also more costly in terms of die space, so that goes at the expense of other things. Also, more complex schemes tend to be slower, which applies here: Intel has the faster cache. Also: beyond 16-wide we get diminishing returns, and 48-way (even 64-way was attempted!) is over the top.

