| | |
|---|---|
| **Stack Machines** | Operand stack to replace accumulator<br>More intermediate results within the CPU -> less memory access required<br>Code density, can reuse to implement recursion<br>Fits naturally to equations (reverse Polish notation) |
| **Reg. Machines** | Small, fast local memory for intermediate results<br>Practical for small memories to be multiported (allows for parallelism in CPU) |
| **Amdahl's Law** | $\text{speedup} = \frac{\text{entire task perf w/ enchance.}}{\text{entire task perf w/out enhance.}}$<br>Make common case fast |
| **CISC** | Eliminate the semantic gap<br>Instruction usefulness limited or so general it is slow |
| **Exceptions** | Hardware (e.g. divide by 0)<br>Software (e.g. SWI #)<br>Invoke exception hander |
| **Memory** | SRAM: maintains store when provided with power<br>4-6 transistors per bit, fast<br>DRAM: requires refreshing<br>1 transistor per bit, fairly fast |
| **Latencies** | Register file: 1 cycle<br>L1 cache: 1-3 cycles<br>L2 cache: 3-9 cycles<br>Main memory: 10-100 cycles<br>Note that DRAM can do burst reads and writes: 2-8 cycles<br>Hard disk: $10^6$ cycles! |
| **Locality** | Temporal: access recently used memory again soon<br>Spatial: access close to recently used memory soon |
| **Cache Line** | Group of around 4 neighbouring words stored |
| **Fully Associative** | Store any word anywhere in the cache, lookup by address<br>Huge overhead involved! |
| **Direct Mapping** | Use part of the address to map onto a cache line<br>Compare tag / valid flag in line to determine if it really caches the address<br>Behaves badly when using data from overlapping addrs. |
| **Set Associative** | As direct mapping, but with a set of cache lines at each location |
| **Victim** | As direct mapped but with a one line buffer to store the last line overwritten |
| **Cache Line Replacement** | Only applied to set/fully associative caches, clearly<br>LRU: requires more info.<br>NLU: pass the "potato" on if a cache line is accessed<br>Random: simple, works well |
| **Writes** | If data is already in the cache then write over it<br>If not in cache then either fetch on write / write around |
| **Write Through** | Data is written to both cache and lower level memory<br>Common on multiprocessor systems for cache coherency<br>Use bus snooping on SMP (does not scale > 2 CPUs) |
| **Write Back** | Data is written to cache only<br>Written to main memory upon replacement, can use dirty bit to prevent this when it is unnecessary |
| **Write Buffer** | Avoid CPU stalling by buffering writes<br>Also perform write merging: can take advantage of burst |
| **Virtual and Physically Addressed** | V + P addresses only differ in upper bits, if cache is no bigger than a page then can use lower bits of V to access the cache without conflict<br>Otherwise use P address for cache or access with V and use P (concurrently found) to compare tags in cache<br>Still have problems with aliasing of P addresses in V |
| **TLB** | Cache recent translations<br>Fully associative cache<br>TLB miss must be looked up<br>Potential control hazard |
| **Pipelining** | Increases frequency and latency (due to latch time)<br>Best to make pipeline stages of similar length |
| **Exceptions** | Imprecise: deep pipeline<br>Precise: when you want OS to be able to restart it |
| **Instruction Replays** | Only find out about cache misses in next stage! |

| | |
|---|---|
| | This means that the pipeline has advanced too far: must replay instructions to give cache hardware more time<br>While waiting for cache, refill the pipeline with instructions from the I-cache |
| **Control Hazard** | Caused by branching and not clearing the pipeline<br>Could document behaviour and use branch delay slots |
| **Data Hazard** | Caused by not taking account of results yet to be written back that is still in other pipeline stages<br>Use feedforward/bypass<br>Sometimes have to stall pipeline if later instructions depend on e.g. memory fetches: can introduce bubbles, so minimise this<br>Could document behaviour and use load delay slots<br>Otherwise have hardware detect it (e.g. scoreboarding) |
| **Parallel** | Synchronous communication<br>Timing assumption can be violated by skew: doesn't work for high f or long s |
| **Serial** | Asynchronous transmission<br>Use 8B/10B coding to guarantee the run length (allows clock recovery), DC balance (allows AC coupling) |
| **Control Flow** | Concurrency simulated via interrupts / scheduler<br>Has to throw away register file, disrupt caching/pipeline<br>Load operations cause stall |
| **Data Flow** | Model dependency graph in CPU, execute it concurrently<br>Inherently concurrent and latency tolerant<br>Easy to take advantage SMP<br>Too much concurrency, so assignment is a problem<br>Makes I/O difficult<br>Ineffective use of very local storage (e.g. register, stack)<br>Scheduling policies simple |

**More information may be needed, check what the exam questions ask on this** ☺