# Caches and the Memory Hierarchy

# Latency, size and price of computer memory
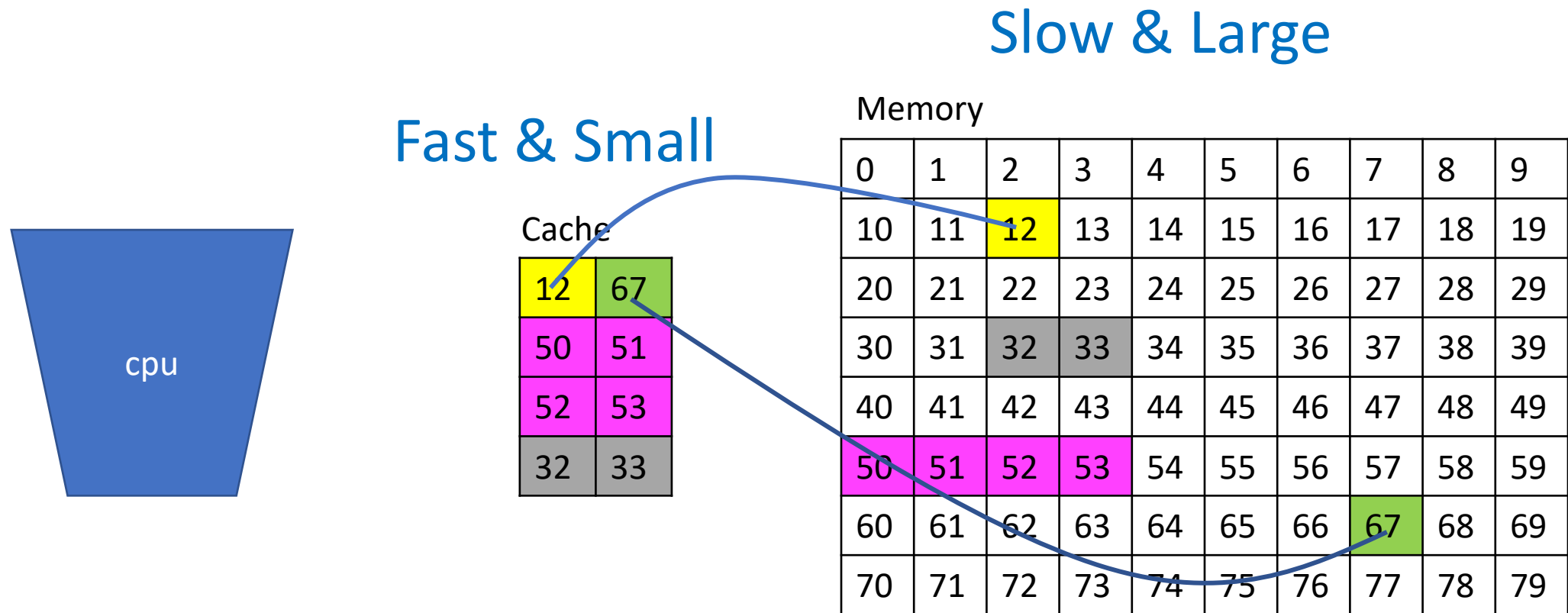
## Given a budget, we need to trade off

$10: Fast & Small
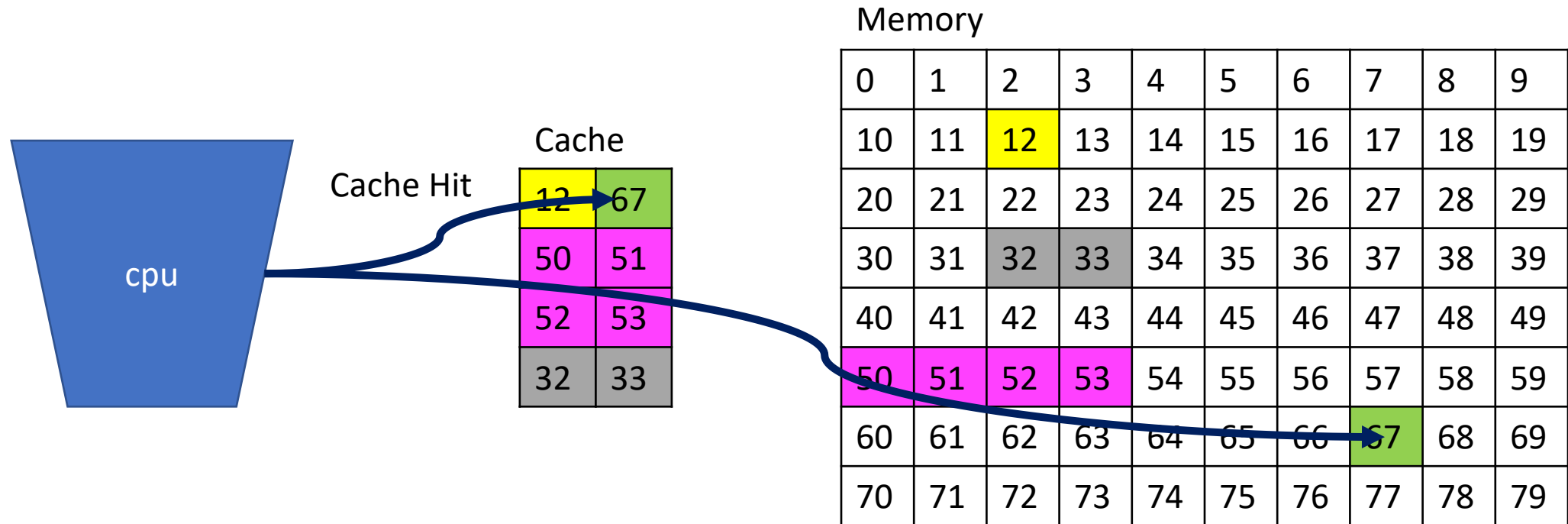
$10: Slow & Large

# Cache: The basic idea
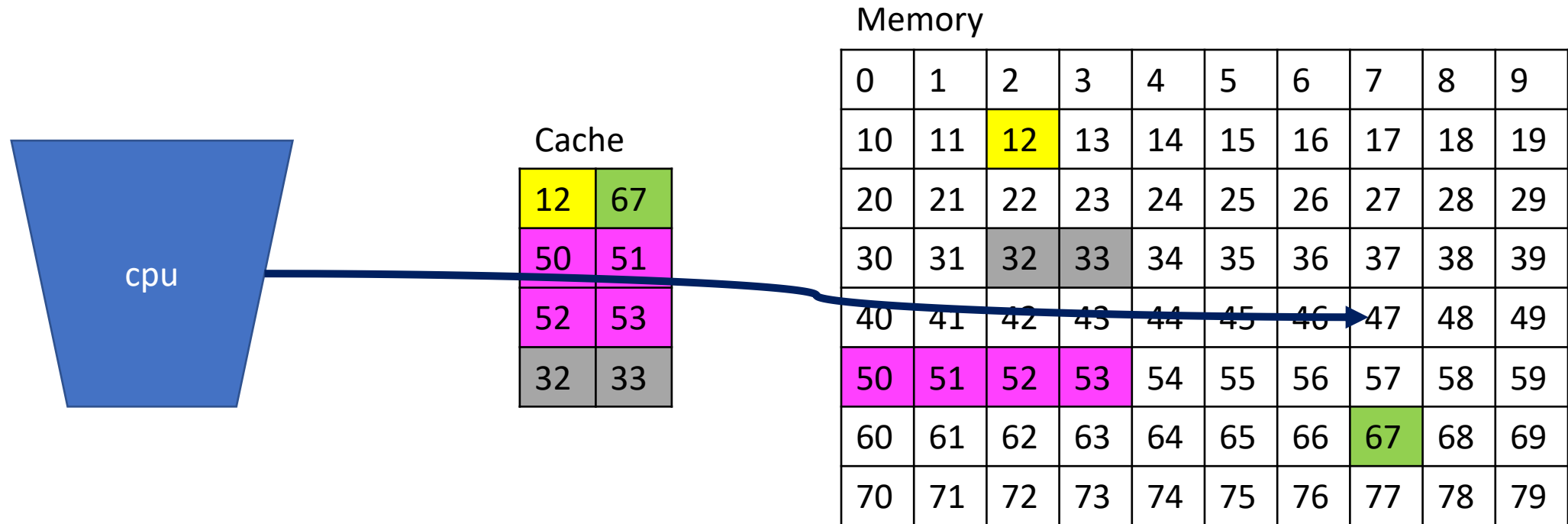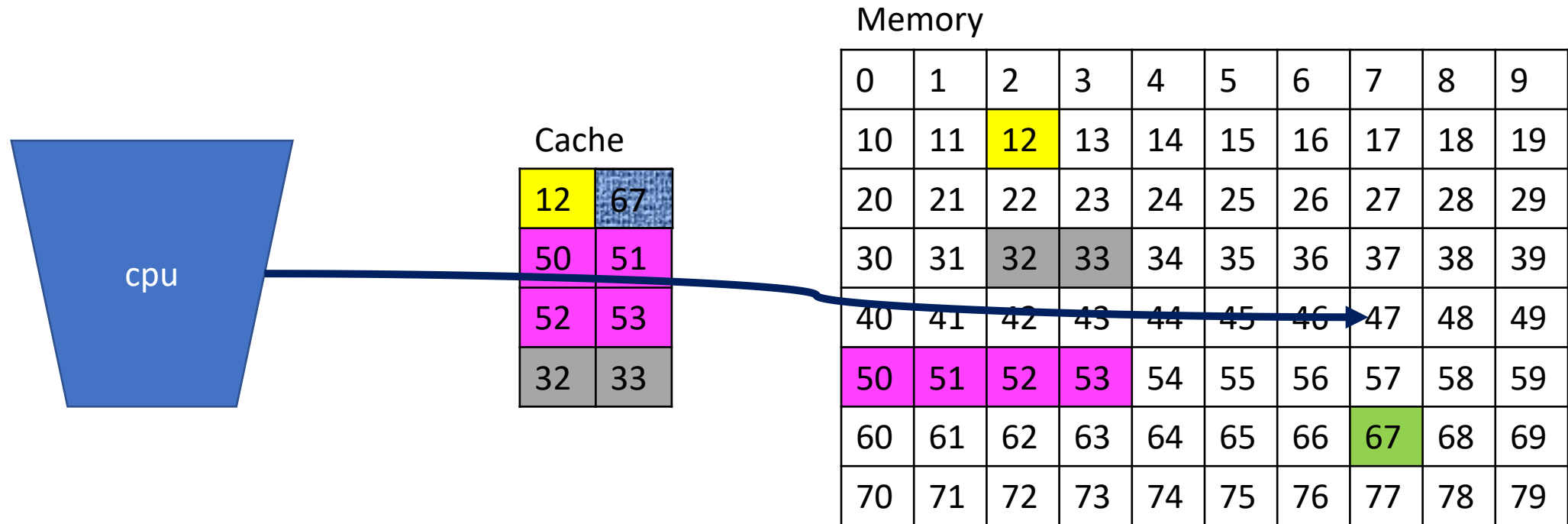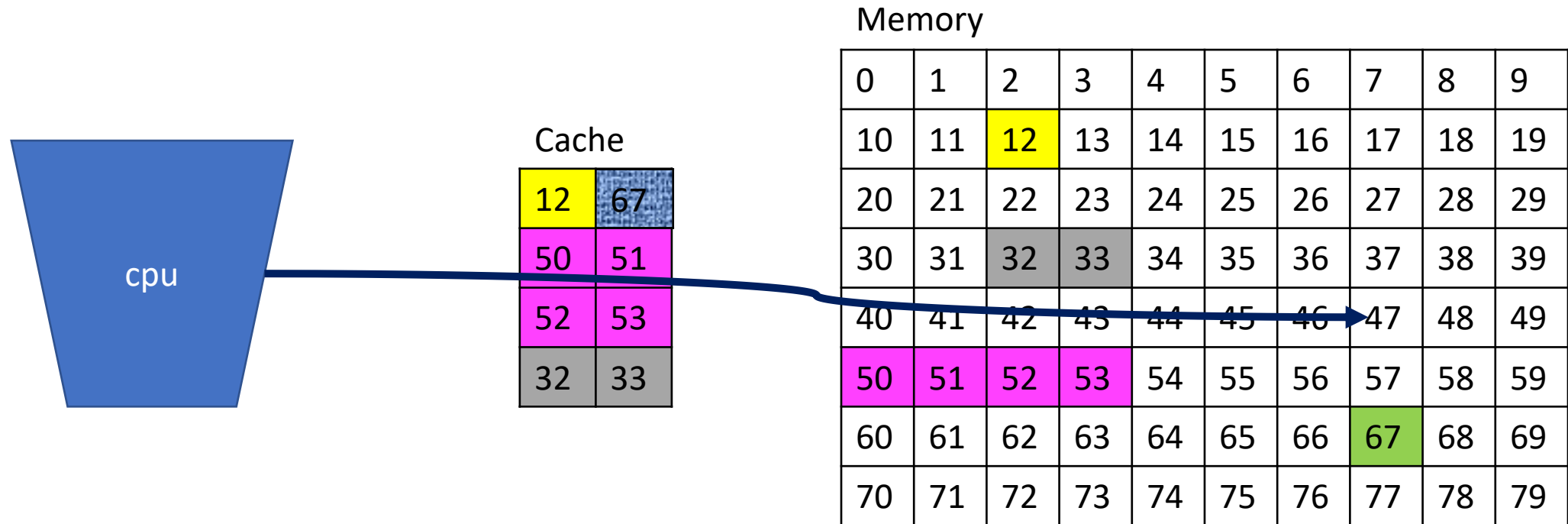


Slow & Large

Fast & Small
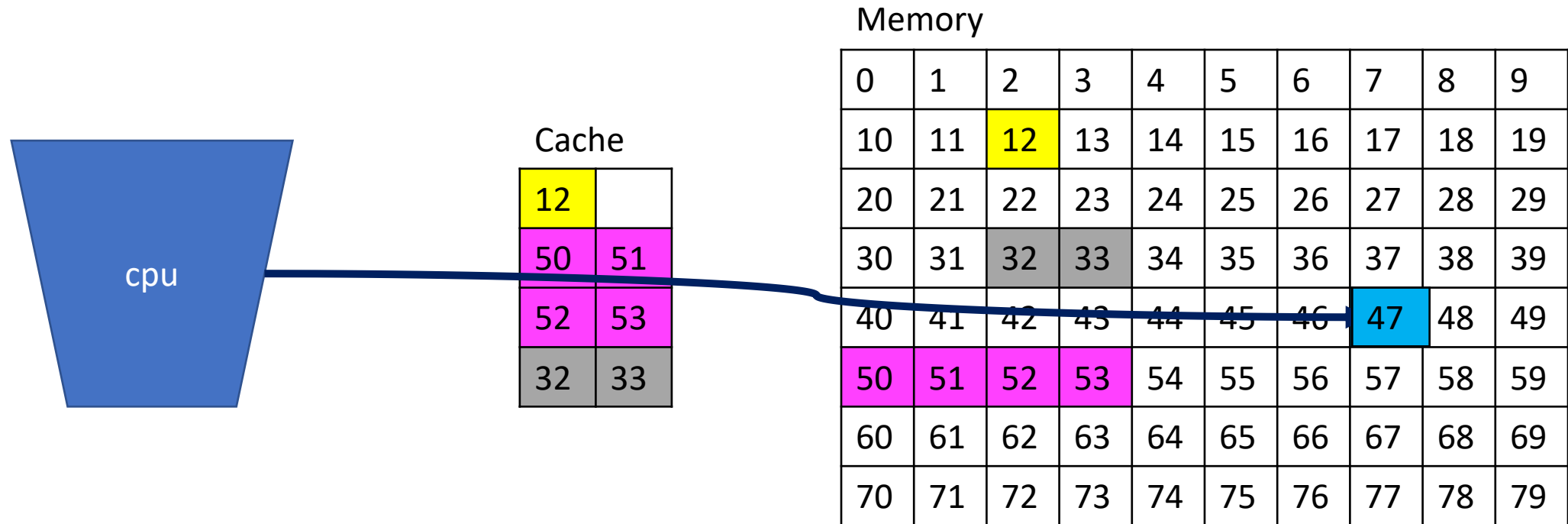
# Cache Hit

# Cache Miss

# Cache Miss Service: 1) Choose byte to drop

# Cache Miss Service: 2) write back

# Cache Miss Service: 3) Read In

# Access Locality

- The cache is effective If most accesses are hits.
    - Cache Hit Rate is high.
- **Temporal Locality**: Multiple accesses to **same** address within a short time period

# Spatial locality

- **Spatial Locality**: Multiple accesses to close-together addresses in short time period.
  - The difference between two sums.
  - Counting words by sorting
- Benefiting from spatial locality
  - Memory is partitioned into **Blocks/Lines** rather than single bytes.
  - Moving a block of memory takes much less time than moving each byte individually.
  - Memory locations that are close to each other are likely to fall in the same block.
  - Resulting in more cache hits.

# Unsorted word count / poor locality

```
=== unsorted list:
the,vernacular,but,as,for,you,ye,carrion,rogues,turning,to,
```

- Consider the memory access to the dictionary D:
- Count without sort:
  D[the]=12332,…,D[but]=943,………,D[vernacular]=10,………….,D[for]=..
- Temporal locality for very common words like "the"
- No spatial locality

# sorted word count / good locality

```
=== sorted list:
lines,lingered,lingered,lingered,lingered,lingered,lingerin
g,lingering,lingering,lingering,lingering,lingering,lingeri
ng,lingering,lingers,lingo,lingo,lining,link,link,linked,li
nked,linked,linked,links,links
```

Entries to D are added one at a time.

1. D[lines]=33

2. D[lines]=33, D[lingered]=5

3. D[lines]=33, D[lingered]=5, D[lingering]=8

Assuming new entries are added at the end, this gives spatial locality.

Spatial locality makes code run much faster (X300)

# Summary

- Caching reduces storage latency by bringing relevant data close to the CPU.

- This requires that code exhibits access locality:
  - Temporal locality: Accessing the same location multiple times
  - Spatial locality: Accessing neighboring locations.