# Four Key Principles from Transformers

## 1. Sequential Processing Not Always Necessary
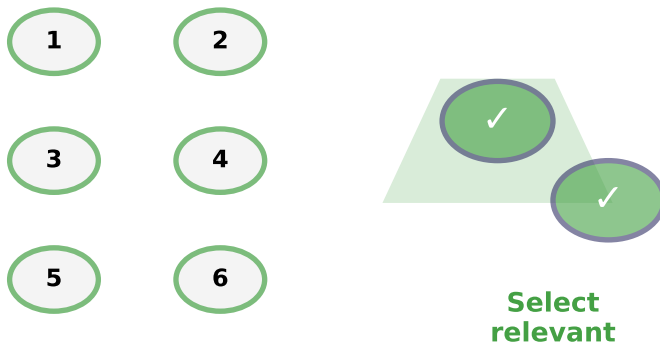
1 → W1
2 → W2
3 → W3
4 → W4
5 → W5

**Order can be encoded,**
not just processed sequentially

## 2. Parallelization Through Independence

P1 → P4
P2 → P5
P3 → P6

**Trade more compute operations**
for less wall-clock time

## 3. Selective Attention vs Compression

1  2
3  4
5  6

✓ ✓

**Select relevant**

**Keep all information**

**Keep information, let model**
decide what's relevant

## 4. Hardware-Algorithm Co-Design

**GPU**
Parallel Processing

✓

**Perfect Match**

**Algorithm**
Transformer Architecture

**Match architecture to hardware capabilities**

Transformer perfectly utilizes
GPU parallel processing power