

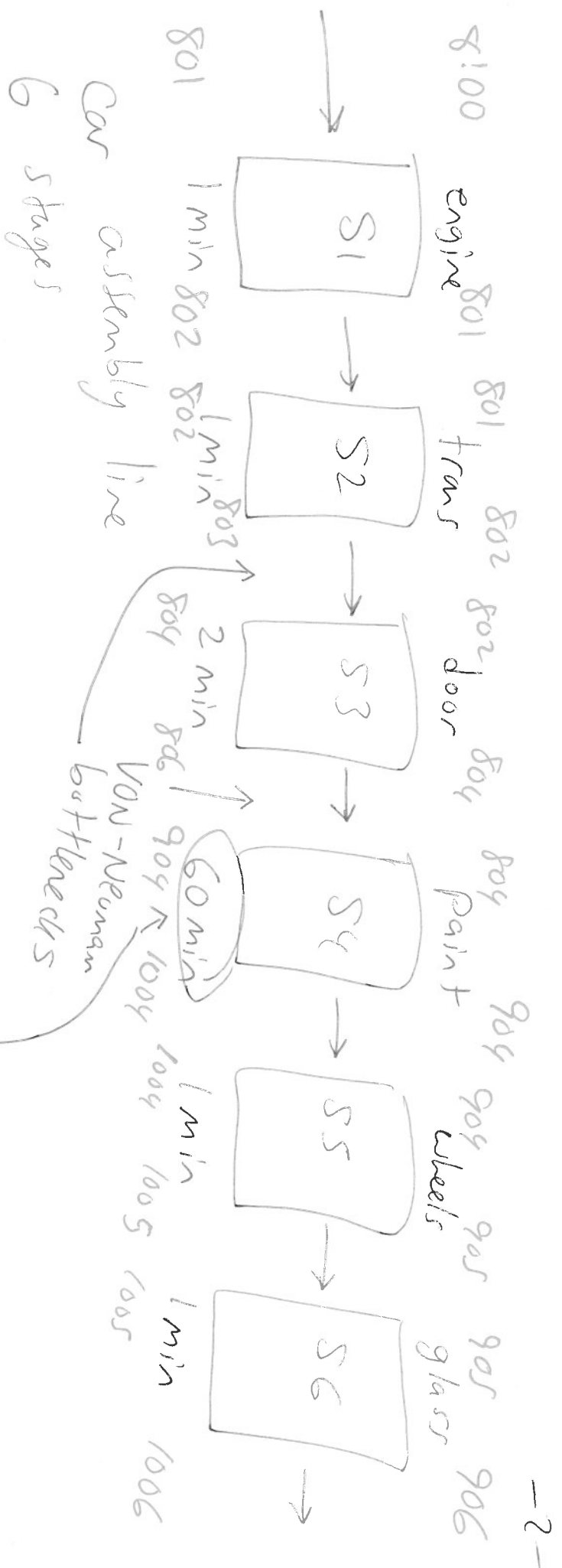
Pipeline:  $\approx$  Assembly line

Multiple hardware working in parallel (and series)  
to fetch, decode, and execute instructions in  
order to maximize bandwidth.

eg MIPS

Bandwidth: # of instructions that are completed  
in a unit of time

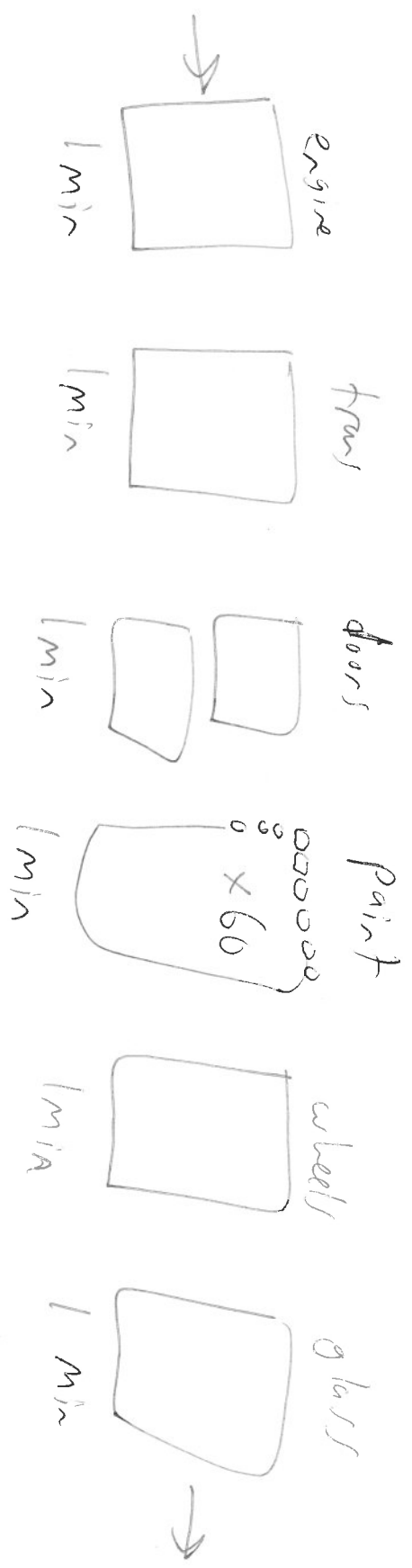
Latency: the total time for one instruction  
to pass through the pipeline from  
start to finish.



Latency:  $\frac{66 \text{ minutes}}{\text{Car}}$

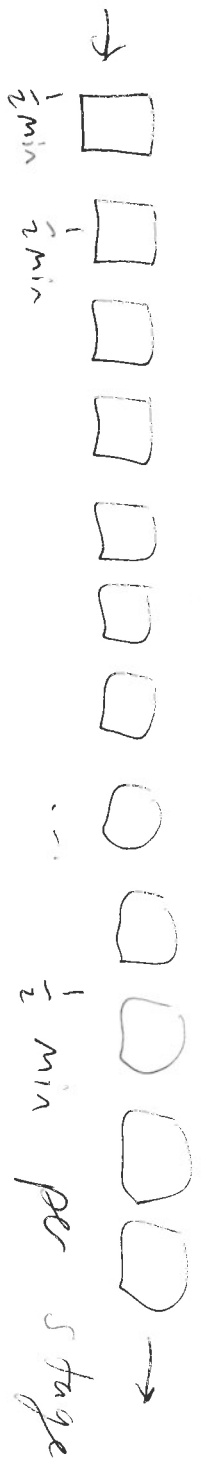
Bandwidth:  $\frac{1 \text{ Car}}{60 \text{ mins}}$

Fix bottlenecks:  
 put extra at slowest stages:



Latency:  $\frac{6 \text{ min}}{\text{car}}$

BW:  $\frac{1 \text{ car}}{\text{min}}$



ie doubled the pipeline length

Latency:  $6 \frac{\text{min}}{\text{car}}$  } doubled the BW  
BW:  $2 \frac{\text{car}}{\text{min}}$  } but same latency  
we like deep pipelines

Disneyland It's a small world 15 min short line

Space Mountain

Latency: want long ride short line  
BW: more riders per hour  
No line } lots of hardware  
long ride }

Oct 5 1995

Jan 30 1937

① 7

② 11

③ 2

④ 5

⑤ 1

① # of 12's in the last 2 digits of the year

② left over

③ # of 4's in left over

④ Add the day

⑤ Add the month

code:

JFM	144
AMJ	025
JAS	036
OND	146

35

% 7 =

0

Sa

26

% 7 =

5

↑

Su	M	Tu	We	Th	F	Sa
1	2	3	4	5	6	0

for 2000's

+1

April 1 1984 Sunday  
Dec 25 1950 Monday } This type of question is on quiz (bonus)

# Chapter 2 of textbook

Pipeline: multiple HW

Purpose: FDS instructions in parallel to maximize BW

Superscalar arch: @ bottlenecks

BW: prefetch, FDS maximize MIPS

Latency: minimize latency



Netflix:

want high BW  
low latency: long load  
fine the  
ok

-7-  
-7-  
-7-

