

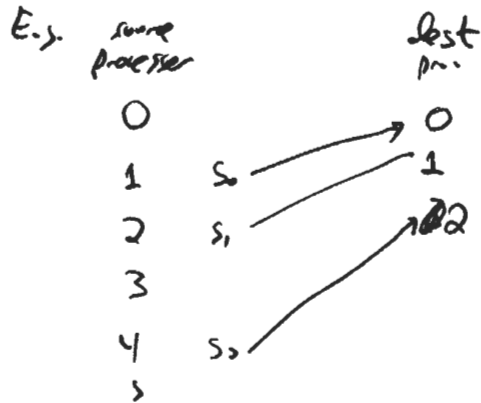
# ONE LAST LECTURE ON ROUTING

I have a magic trick too! (Nerd magic)...

First, I need to define a squish message pattern.

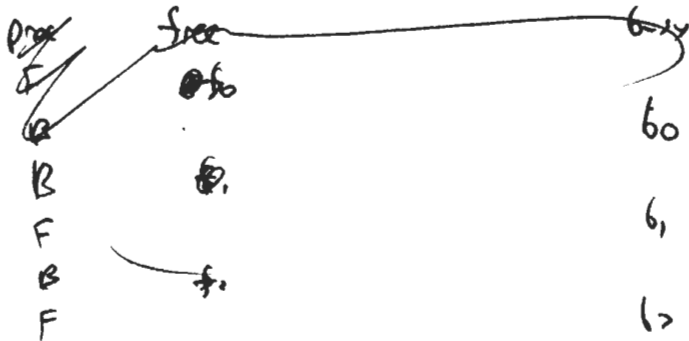
Given a set of processors  $\{s_0, \dots, s_{n-1}\}$   
with  $0 \leq s_0 < s_1 < \dots < s_{n-1} < P$

Processor  $s_i$  sends a message to processor  $i$ .

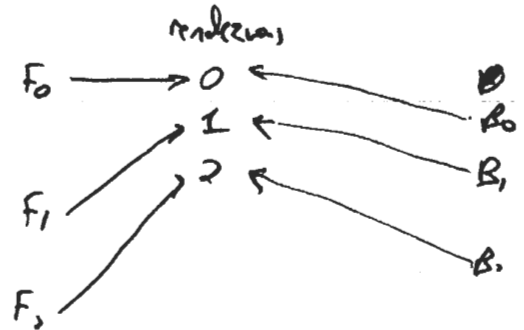


Why would I be interested in squish?

Could use it to pair up free + busy processors



Processor	State
0	F <sub>0</sub>
1	B <sub>0</sub>
2	B <sub>1</sub>
3	F <sub>1</sub>
4	B <sub>2</sub>
5	F <sub>2</sub>



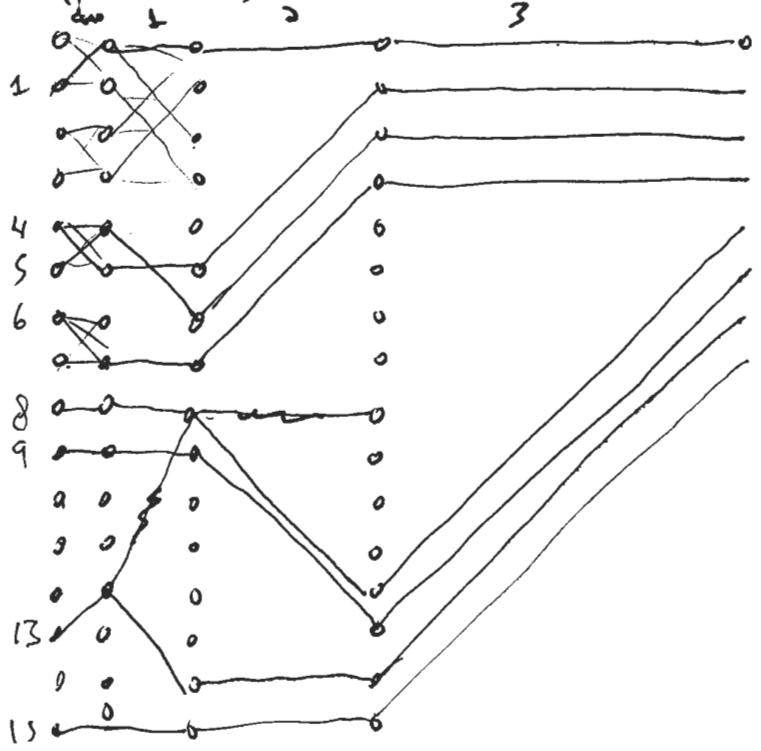
Now processor 1 knows the address of 1 free + 1 busy processor.

Could be used in a cilk implementation for example.  
Or parallel "cars"

Now the magic: I have a butterfly



Now the magic! Here is a butterfly  
(pre draw of bit wire)



Pick some subset that I will

Squash e.g. 1, 4, 5,  
6, 8, 9, 13, 15

$S_0=1$     $S_3=6$     $S_6=13$   
 $S_1=4$     $S_4=8$     $S_7=15$   
 $S_2=5$     $S_5=9$

1101  
0110  
1011

Look MA, NO collisions!

~~This always~~

Then: This magic always works

Proof: Define: a semicontraction is a 1-1 routing  
from  $\{s_0, \dots, s_n\}$  to  $\{d_0, \dots, d_n\}$

with  $s_i \rightarrow d_i$  such that

$$\forall (i, j) \quad |s_i - s_j| \geq |d_i - d_j|$$

↑ this is subtracting, not knowing distance.

Lemma: Semicontraction routes on butterfly w/o conflict.

proof by contradiction:

Suppose  $\exists$  conflict.  $\Rightarrow$  some pair must conflict

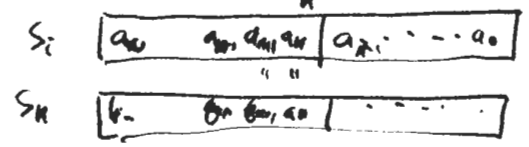
$s_i \rightarrow d_i$  conflicts with  $s_j \rightarrow d_j$

looks like



Observe  $|s_i - s_j| < 2^k$

why?  $s_i + s_j$  must agree on bits high bit



but  $a_i = b_i$  for  $i \geq k$

$$\text{so } |s_i - s_j| = \left| \sum_{i=0}^{k-1} a_i 2^i - \sum_{i=0}^{k-1} b_i 2^i \right|$$

$$= \left| \sum_{i=0}^{k-1} (a_i - b_i) \cdot 2^i + \sum_{i=k}^n (a_i - b_i) 2^i \right|$$

$$= \left| \sum_{i=0}^k (a_i - b_i) \cdot 2^i + 0 \right|$$

$$\leq |2^k - 1| < 2^k$$

observe  $|d_i - d_j| > 2^k$

most agree on low bits

$$d_i = \overbrace{\dots c_{k+1} c_k}^k c_0$$

$$d_j = \overbrace{\dots e_{k+1} e_k}^k c_0$$

$$c_i = e_i \text{ if } i \leq k$$

$$\begin{aligned} |d_i - d_j| &= \left| \sum_{i=0}^k (c_i - e_i) 2^i \right| \\ &= \left| \sum_{i=0}^k (c_i - e_i) 2^i + \sum_{i=k+1}^n (c_i - e_i) 2^i \right| \\ &= \left| 0 + \sum_{i=k+1}^n (c_i - e_i) 2^i \right| \end{aligned}$$

some bit must be different, so  $> 2^k$

$\Rightarrow$  not a semicontraction  $\neq \boxtimes$

Observe that squash is a semicontraction and  
 $\therefore$  The magic works  $\boxtimes$

What can we do with this magic?  
 How do we implement it?

Application of magic: routing 1-1 message traffic

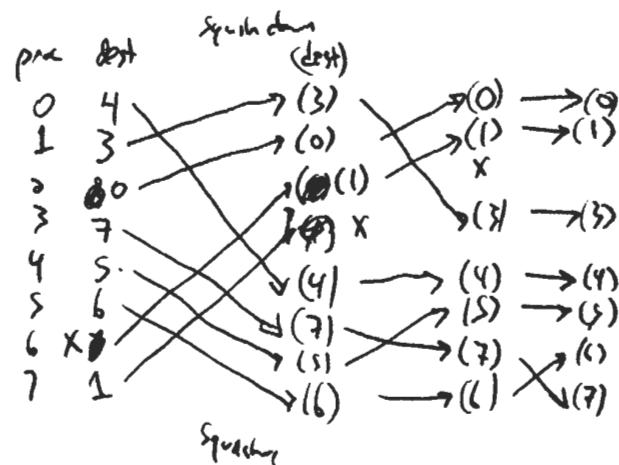
Algorithm:

pick all.

squash all messages with 0 in bit (n-1).

squash up all messes with 1 in bit (n-1)

+ recurse



# How to implement $\text{scan}$ .

The problem is to compute the enumeration of the set  $\{s, \dots\}$

$$\text{let } v_i = \begin{cases} 0 & \text{if } i \notin S \\ 1 & \text{if } i \in S \end{cases}$$

i.e.

0		0	
1	$s_0$	1	
2	$s_1$	1	
3		0	
4	$s_2$	1	
5		0	

also  $\text{scan} +$

Want to compute the prefix sum  $w$

$w_i =$  the number of 1's in  $v$  before  $i$

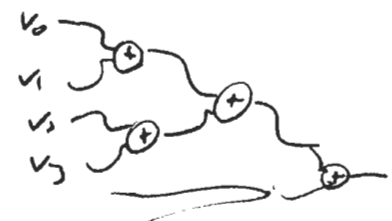
$$w_i = \sum_{j=0}^{i-1} v_j$$

Can compute  $w$  in work  $\Theta$  and  $T_{\text{op}} = \Theta$   
 Can we do it faster?

Yes. For example we can compute ~~it~~

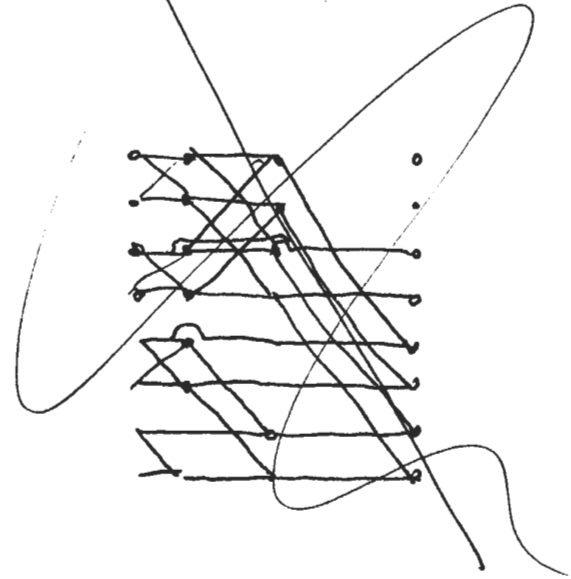
$$w_p = v_0 + v_1 + \dots + v_{p-1}$$

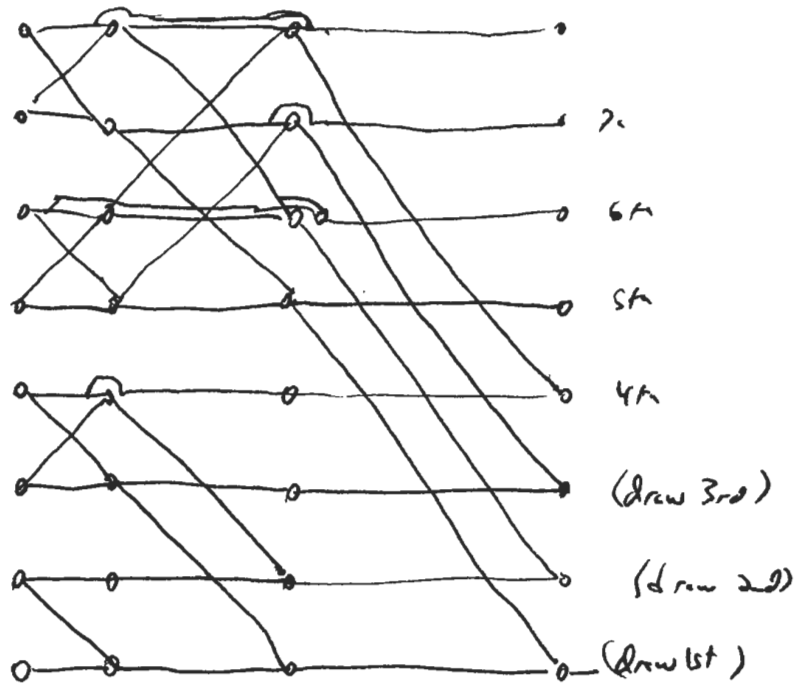
$$= ((v_0 + v_1) + (v_2 + v_3)) + ((v_4 + v_5) + \dots)$$



$T_{\text{op}} = \lg n$   
 but work?

On a hypercube! Butterfly! just do it,  $\leftarrow$  easy  
 common subexpressions



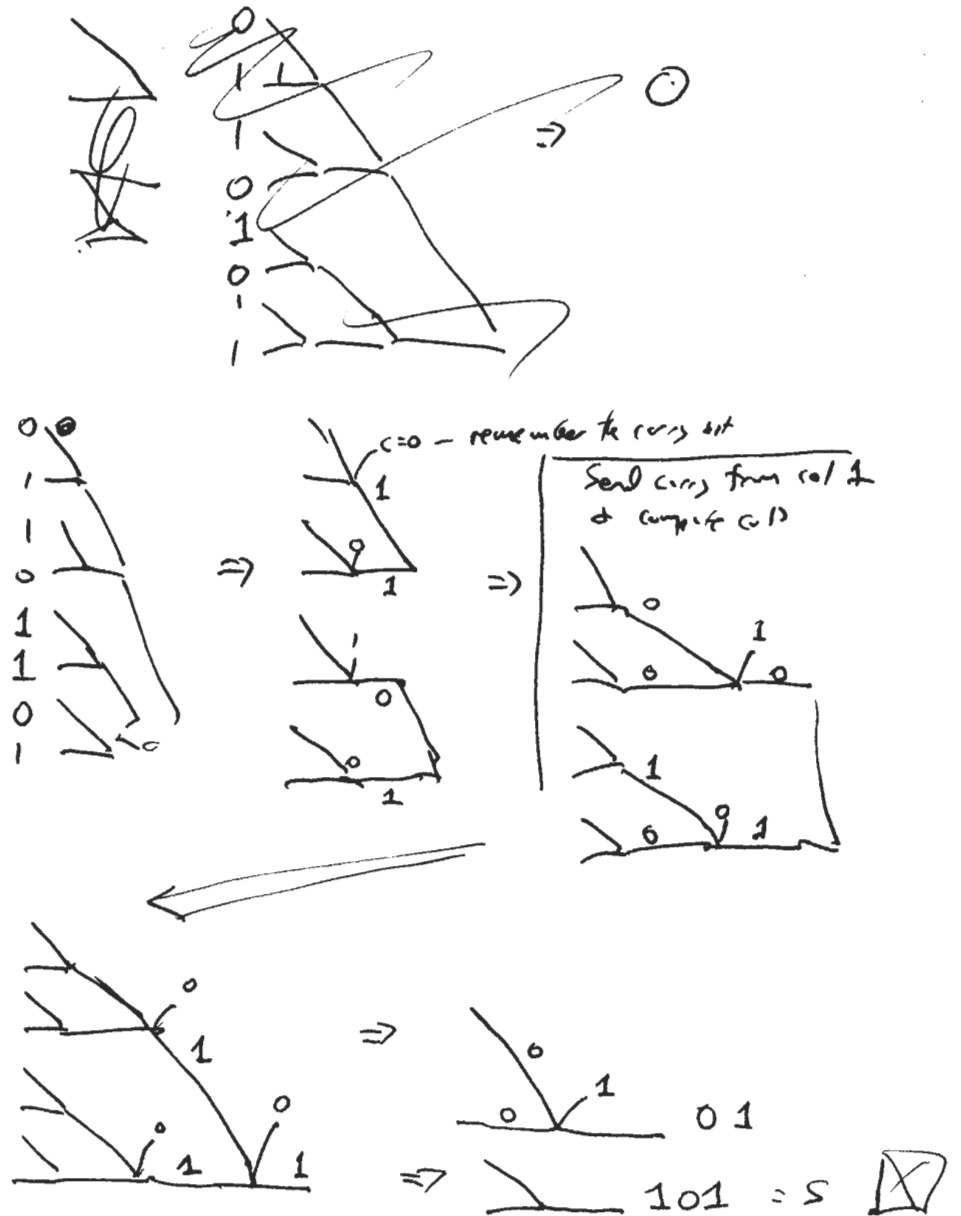


It fits in a butterfly!

Analysis: Time to do scan +  
 depth of graph =  $\lg n$   
 how long to add 2 numbers?  
 we assume the wire is  $O(1)$   
 bits/time unit  
 $\Rightarrow$  it takes  $O(\lg n)$  time  
 just to get the final answer  
 at the last stage.

But, the adds can be pipelined!  
 so is it  $O(\lg^2 n)$  time?

No! Adds can be pipelined.



With pipelining it is  $O(\lg n)$  for / +

~~Total time to run a 1-1 routing~~

$$T_n = O(\lg n) + T_{n/2}$$

do 1-1 routing on  $n$  processors:

compute scan + up on  $O$ 's  $O(\lg n)$

Send message in squish down  $M + \lg n$

compute scan + down on  $O$ 's  $O(\lg n)$

Send message in squish up  $M + \lg n$

do 1-1 routing on two subcubes in parallel

$$T_n = 2M + O(\lg n) + T_{n/2}$$

$$T_n = O(M) + T_{n/2}$$

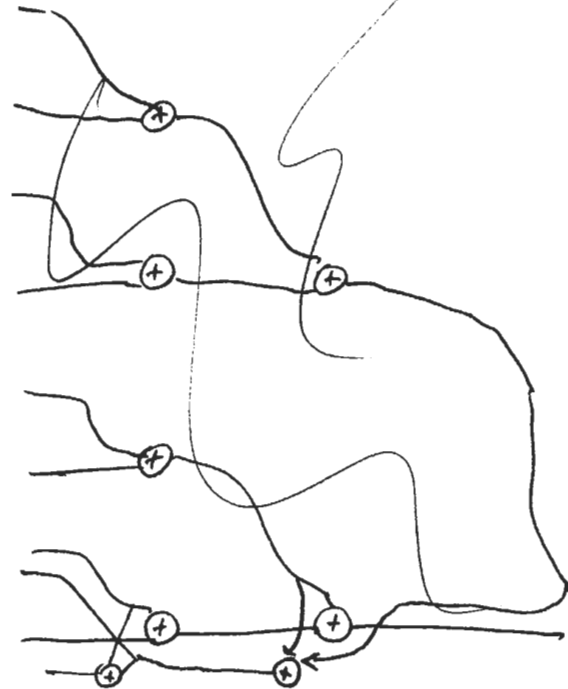
$$= O(M \lg n)$$

$$M = \Omega(\lg n)$$

$$O(M) = \lg n$$

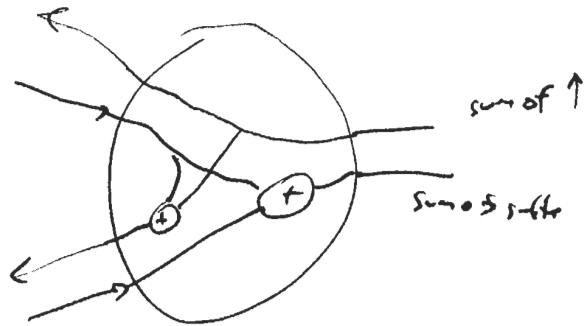
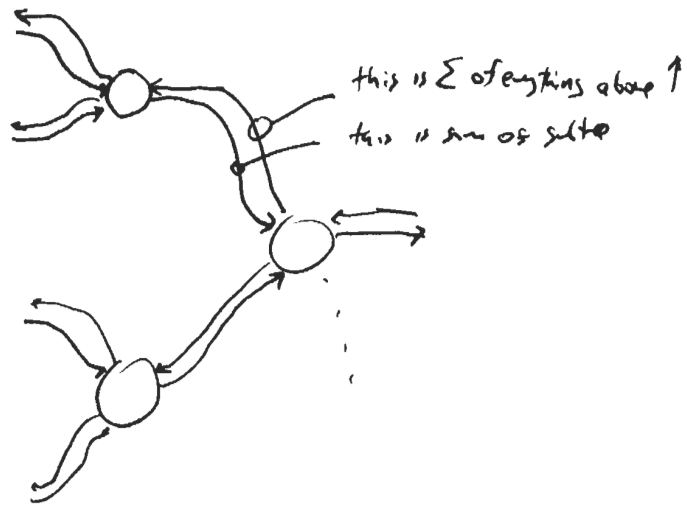
19.6

Another way to do parallel prefix with two trees  
super important



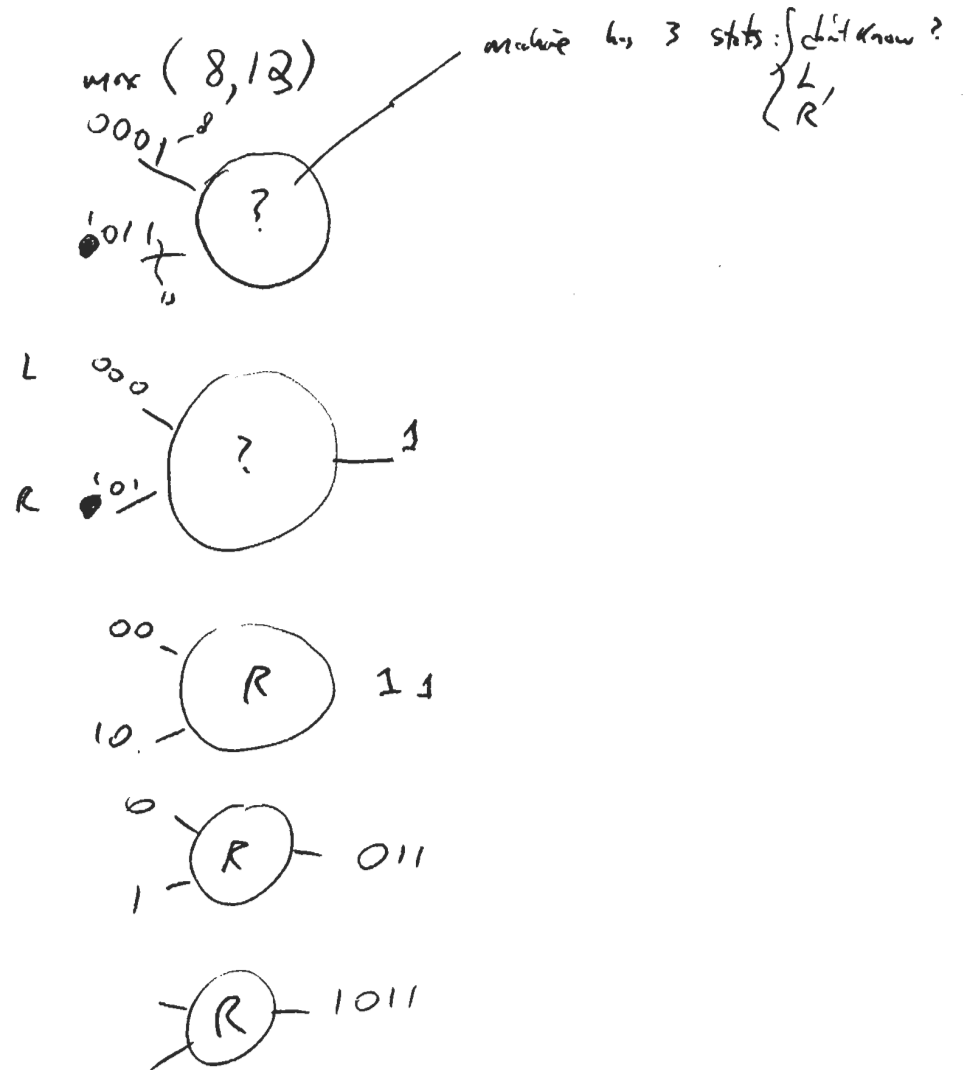
Another way to do  $\uparrow$   
with depth  $2 \log p$  adders

but work only  $p$   
(butterfly is depth  $\log p$ , work  $p \log p$ )



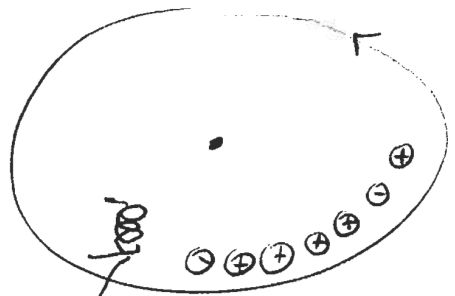
1.1.7

Also can bit-pipeline max; feed in most significant bits first



# Disks:

A spinning platter covered with magnets



encode info in the orientation of magnet.

To write: position an electromagnet over a spot and apply current, it switches the state

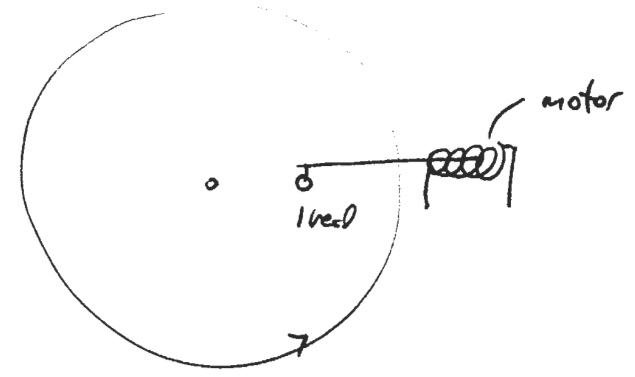
To read: in principle, the ~~same~~ goes in reverse. sense current from the induced field.

## Realities: ~~Signals/Magnets~~

GMR - Giant Magnetoresistance  
a device whose resistance is a function of magnetic fields. Very sensitive.

GMR discovered in late 80's  
=> about 10 years to appear in disks.  
About 20 ~~MB~~ Gb/in<sup>2</sup>

# How to position the head?



disk spins  
motor moves head radially ("seek")  
=> can see any spot on disk

## Motor is a linear induction motor

(old motors were stepper motors. time to seek distance  $n$  was  $\Theta(n)$ )

with linear induction motor, head accelerates to 1/2 any point => time  $\Theta(\sqrt{n})$

today, avg seek time = ~~2.9ms~~ 4.9ms read, 5.4ms write  
best seek time = ~1ms  
worst = ~9ms

rot = 10,000 rpm (expensive + hot) 6ms ↘  
5,400 rpm 11ms ↘  
(1/2 time on average.)



More data near edge of disk

⇒ trick: put your data on the  
outer part of the disk to  
reduce # of seeks.

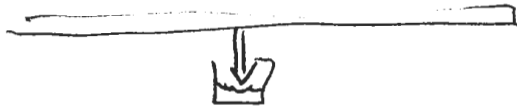
Head "flies" above disk surface



Disk bearings use fluid-dynamics too  
(last ball bearings - about 2 years ago)

- Quiet
- No rumble (vibration, moving around, heat, etc.)  
Non-repetitive rumble

the disk itself "rums" in



with ball bearings heat + vibration  
were the disks big energy consumers  
to bearings to deform.

Now - I don't know what causes heads  
to wear out.

~~Observation to work I remember but is~~

Drive in use

488

458-799 MB/s (up to 200 MB/s)

Since avg seek is  $\frac{5ms}{8ms}$  seek

⇒ 16 MB/s transfers needed to  
keep the disk 1/2 busy!