

# Validation via information theory

# Organization

What is PA?

Rationale of PA

Roadmap

Shannon's coding theorem

Formalization of PA

## Posterior agreement

How to validate algorithms ?

A

:

X

distr.  $p(\cdot | X)$  over  
solutions

clustering

dataset  
of points

cluster assignments

shortest path  
from A to B

graph

paths from A to B

linear regression

dataset  
of points

linear models

## Posterior agreement

How to validate algorithms?

$$\lambda: X \mapsto p(\cdot | x)$$

Expected log posterior agreement:

$$\frac{1}{\log|C|} \mathbb{E}_{X', X''} [\log(G|\kappa(x', x''))]$$

where

$$\kappa(x', x'') = \sum_c p(c|x') p(c|x'')$$

$$= \mathbb{E}_{p(\cdot|x')} [p(\cdot|x'')]$$

## Posterior agreement

Expected log posterior agreement:

$$\frac{1}{\log |C|} \mathbb{E}_{X', X''} [\log (C | \kappa(X', X''))]$$

In practice, emp. log PA =  $\frac{1}{\log |C|} \log (|C| | \kappa(X', X''))$ .

When comparing  $A_1$  and  $A_2$ , choose the one  
that maximizes  $\kappa(X', X'')$ .

# Organization

What is PA?

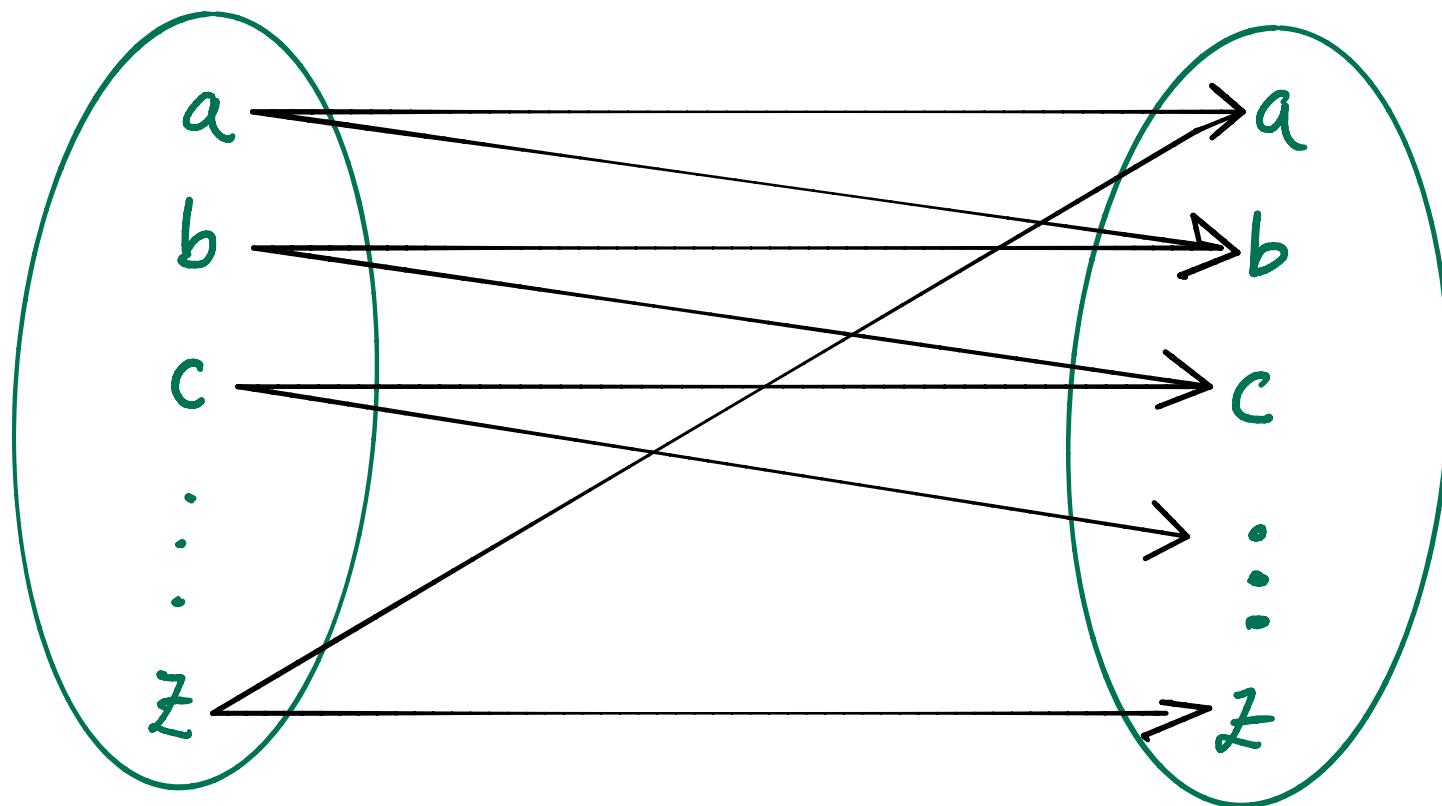
Rationale of PA

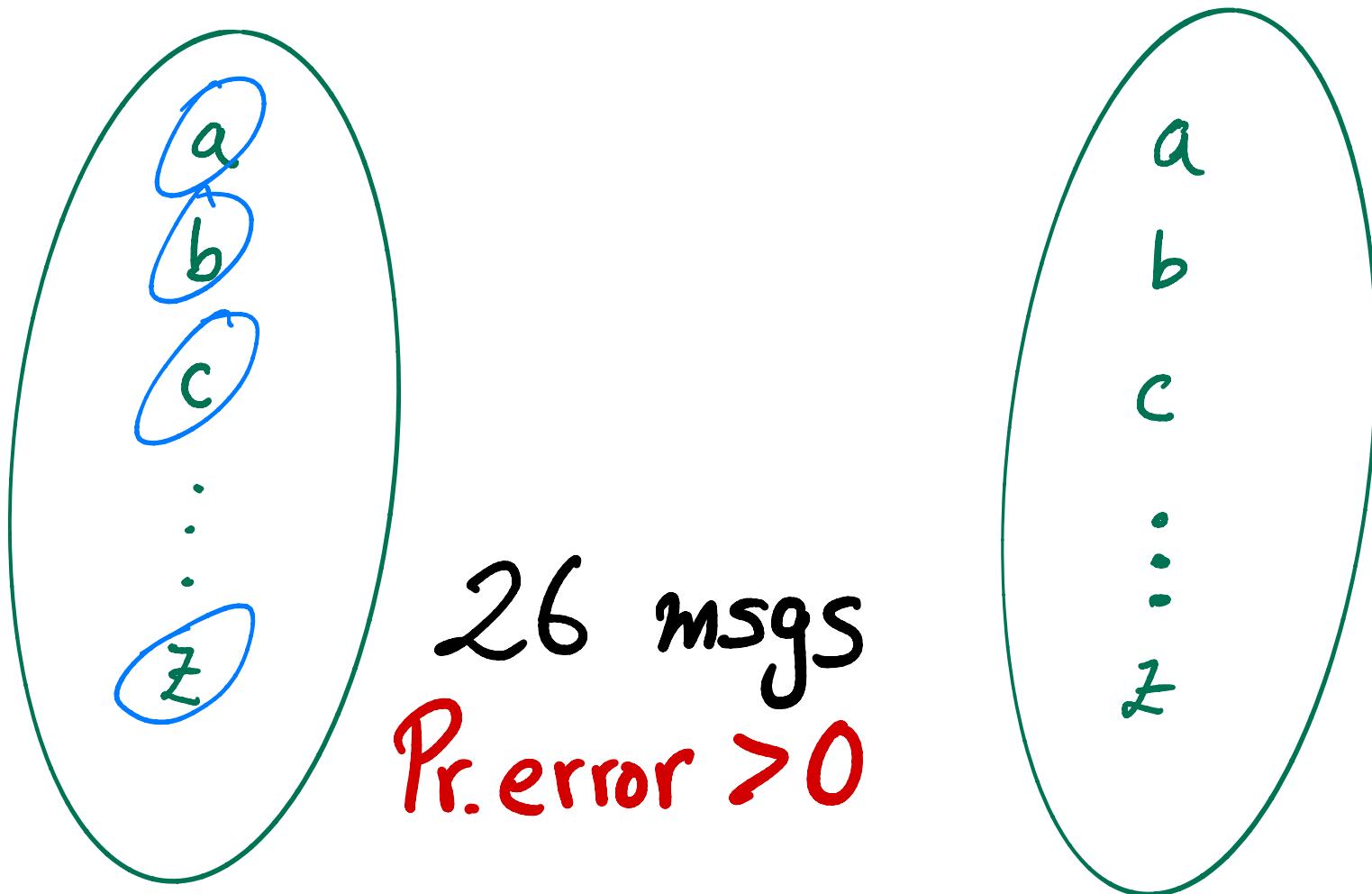
Roadmap

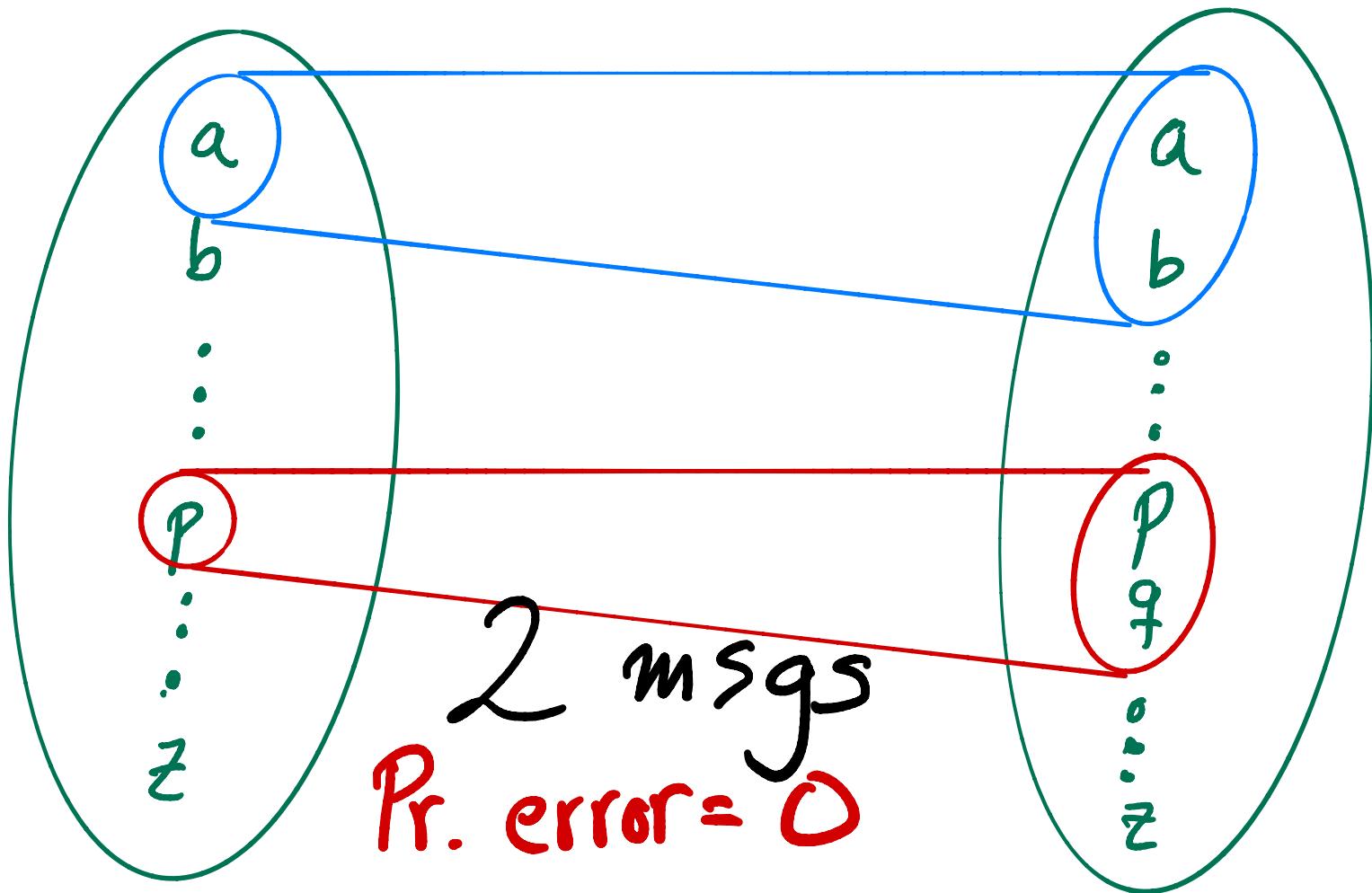
Shannon's coding theorem

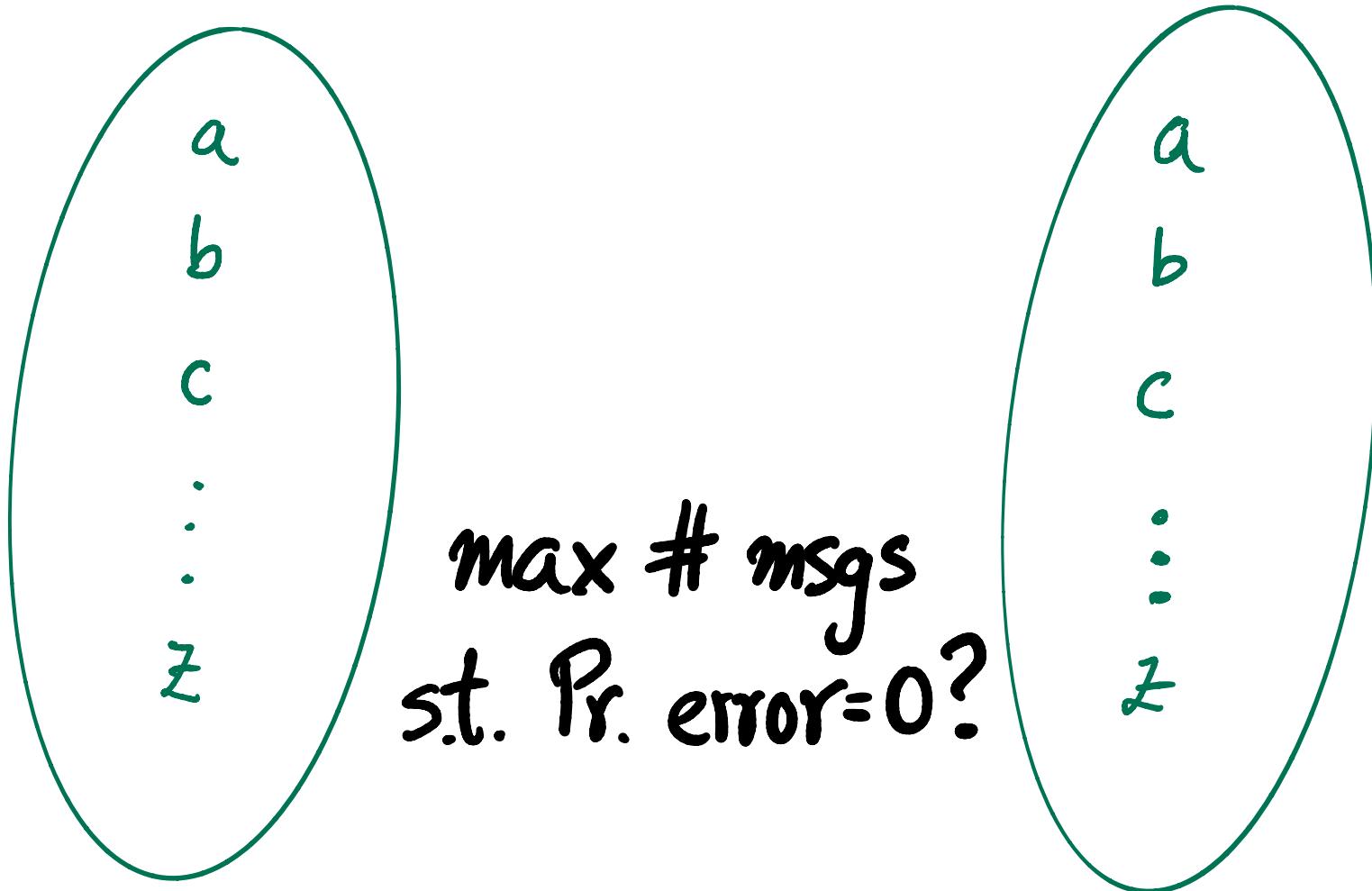
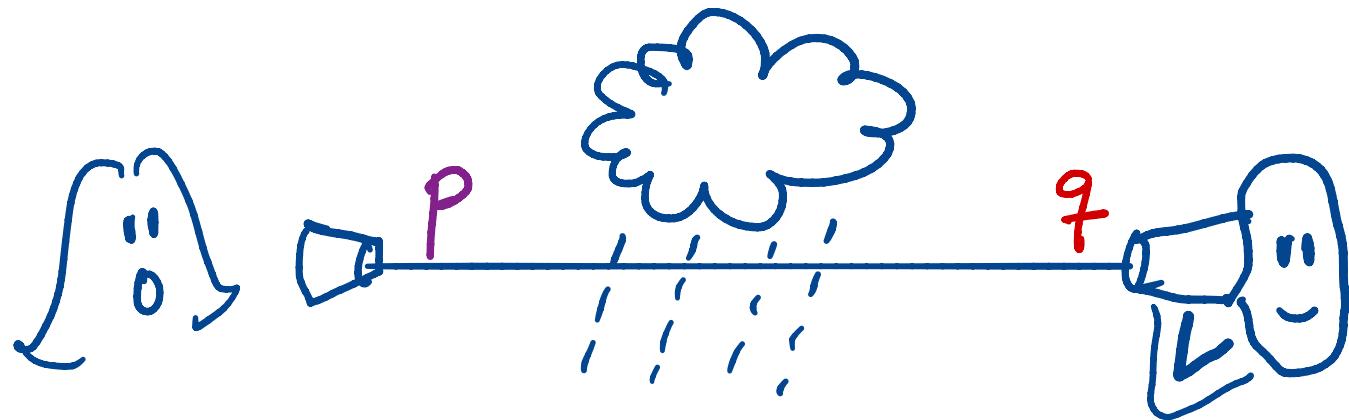
Formalization of PA

# Communication on a noisy channel

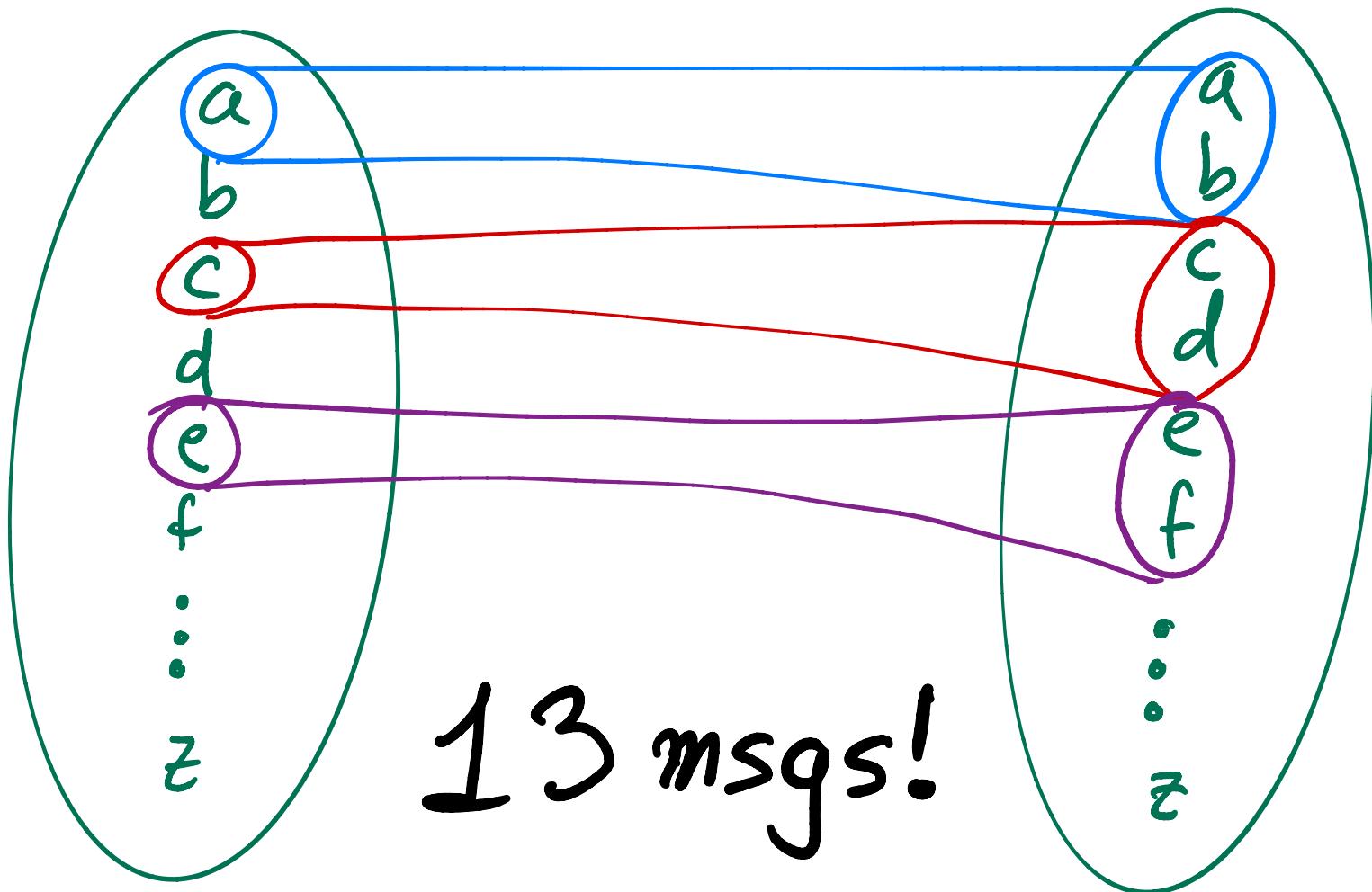
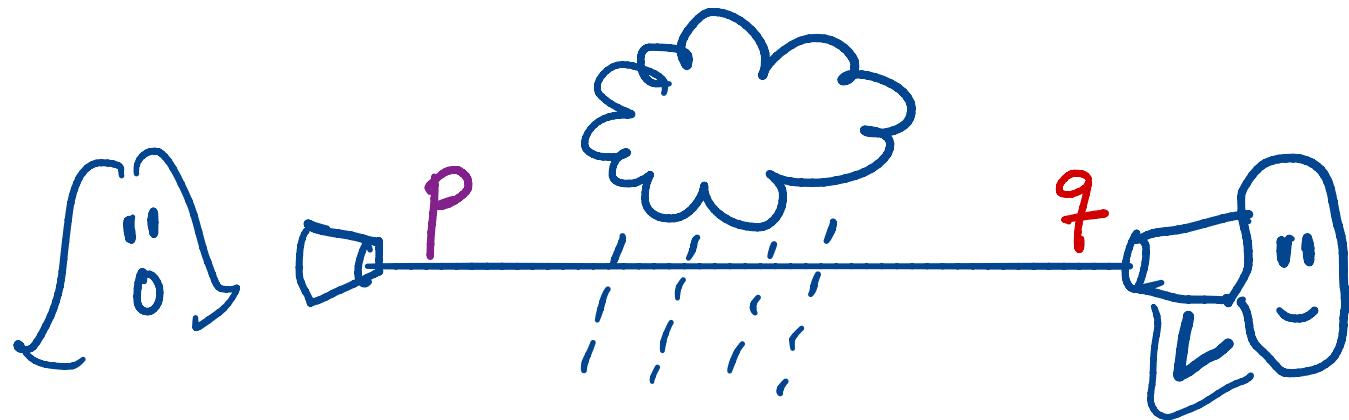


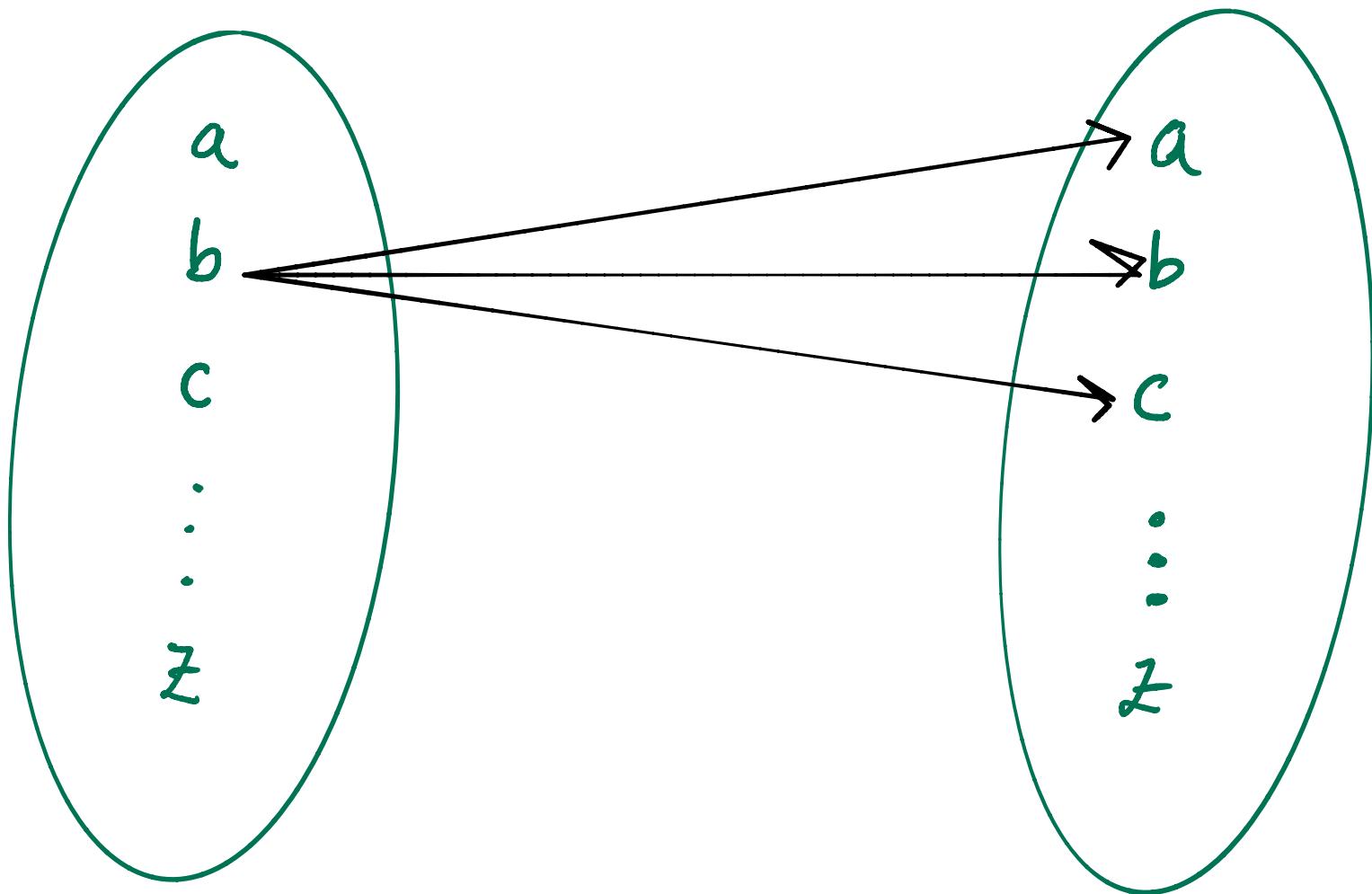
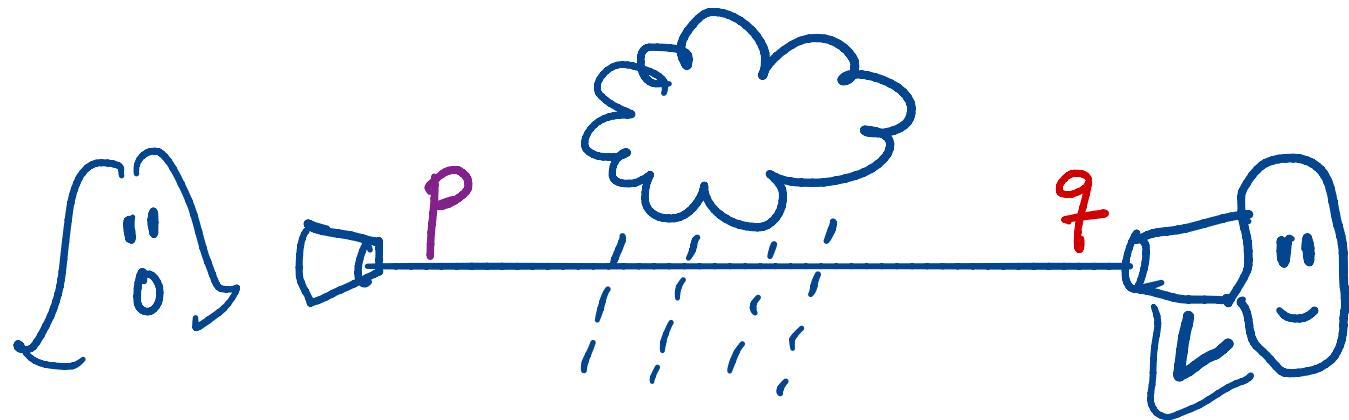


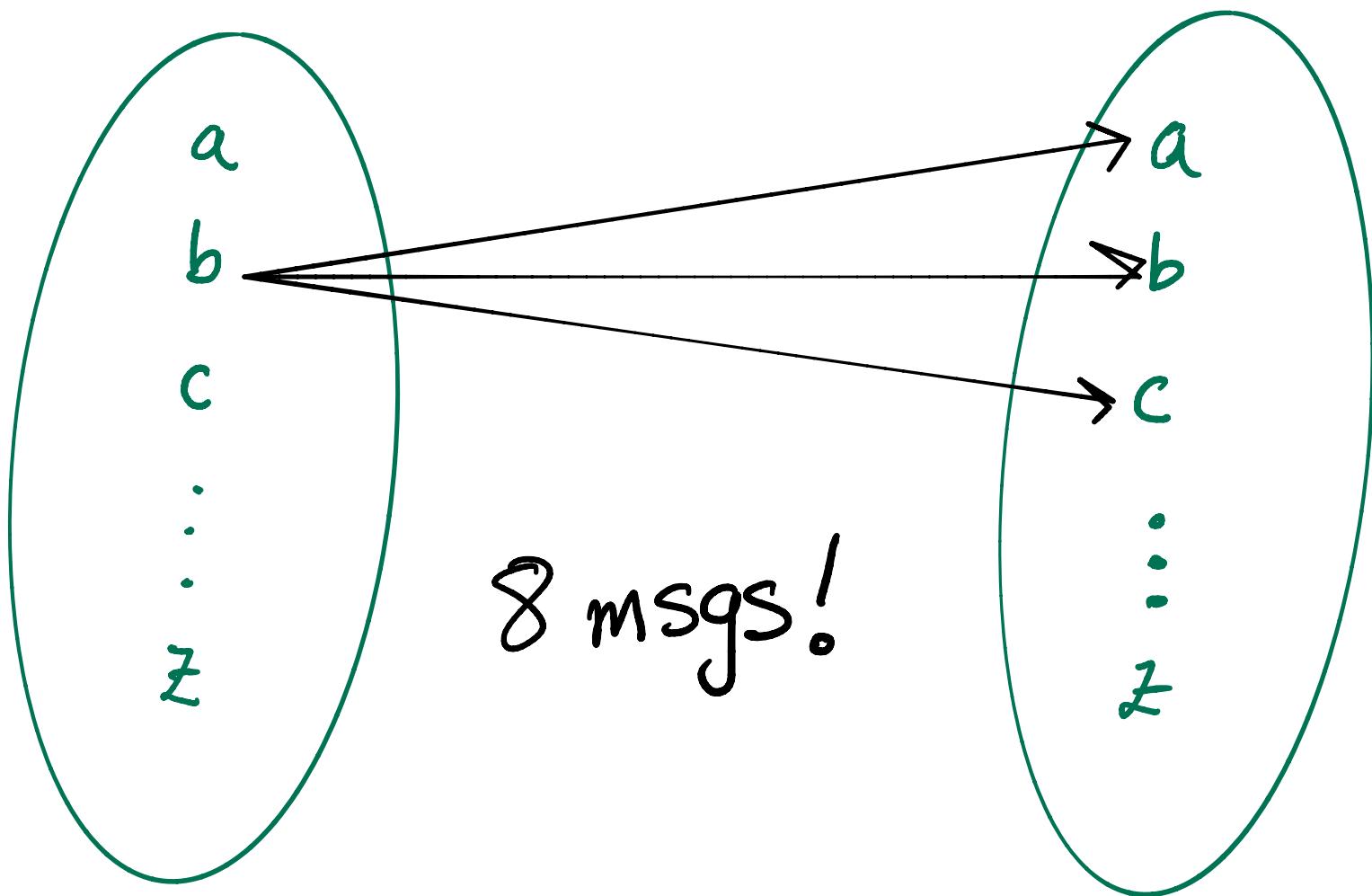


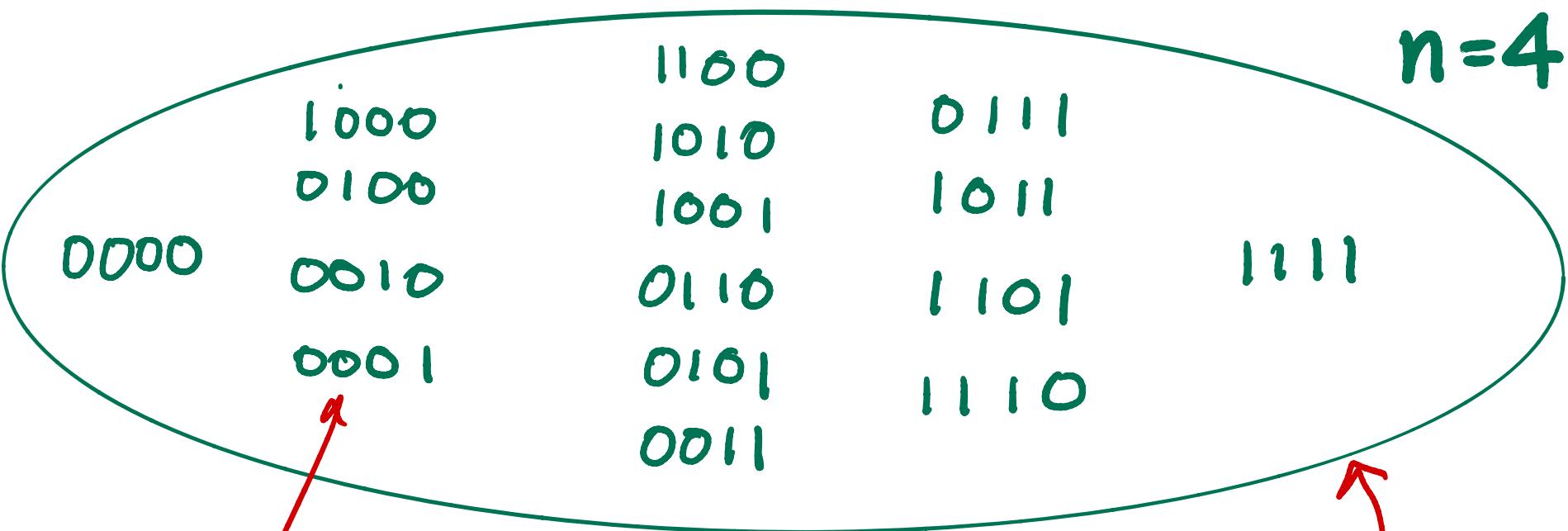


max # msgs  
s.t.  $\Pr.$  error=0?





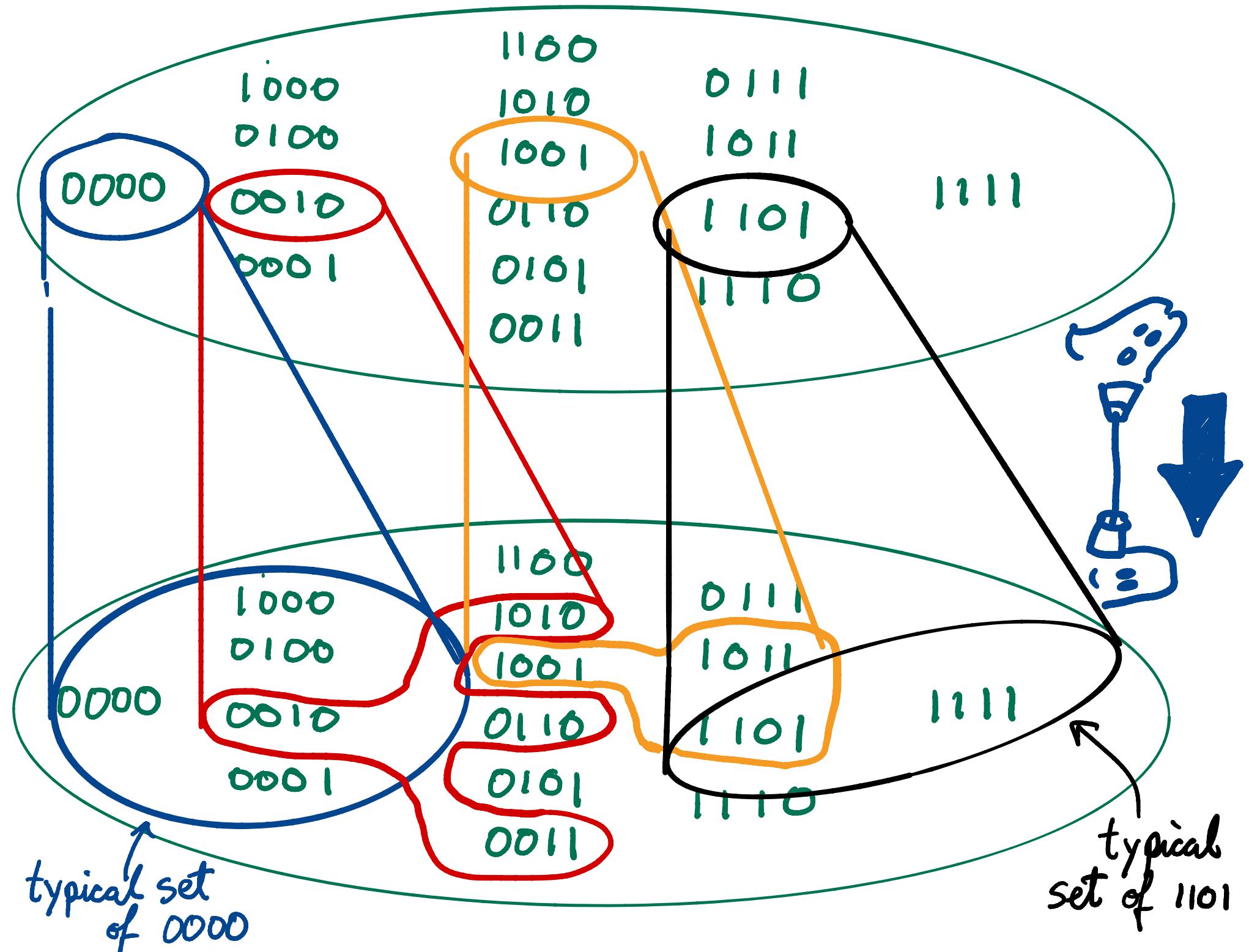


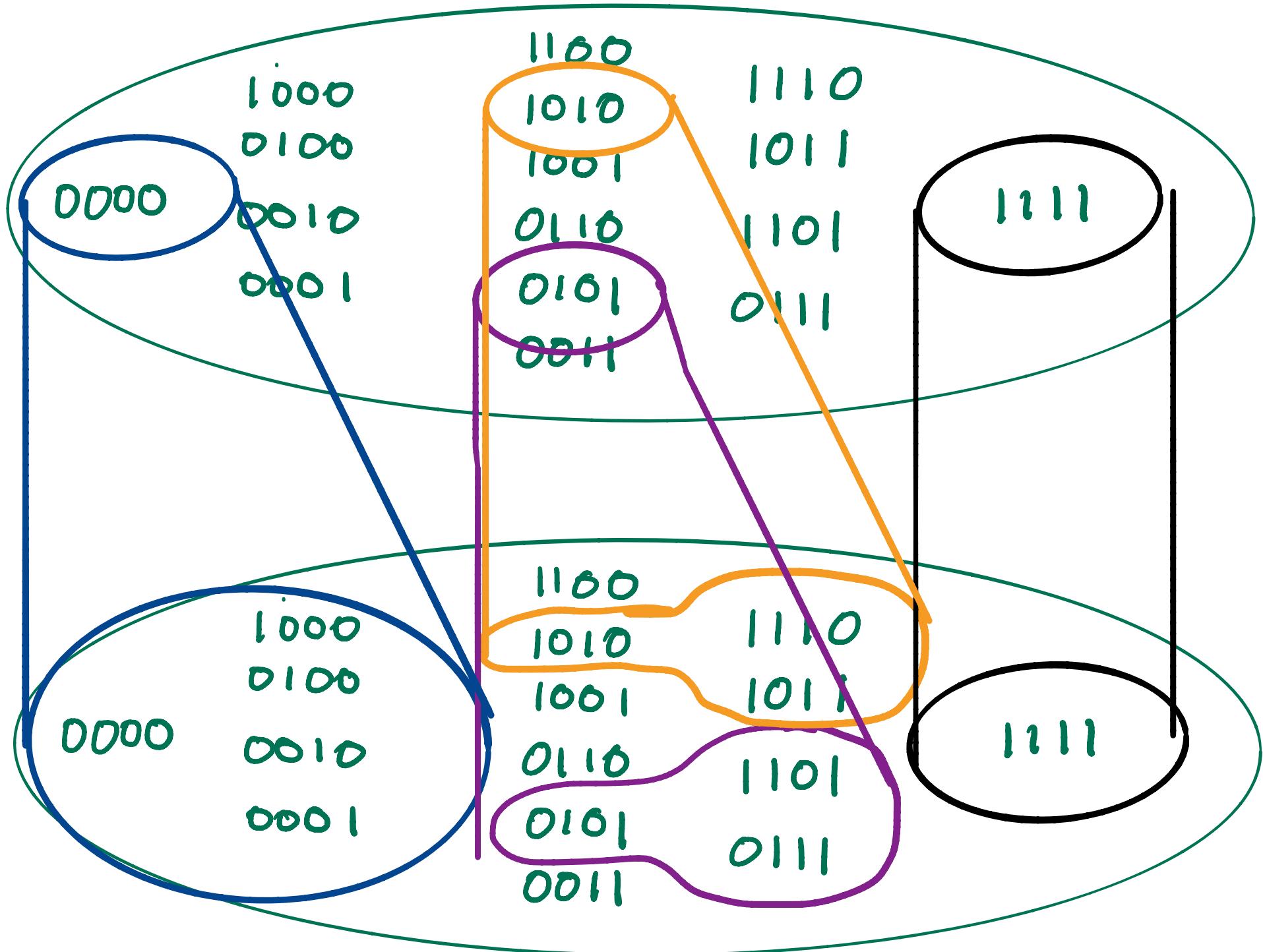


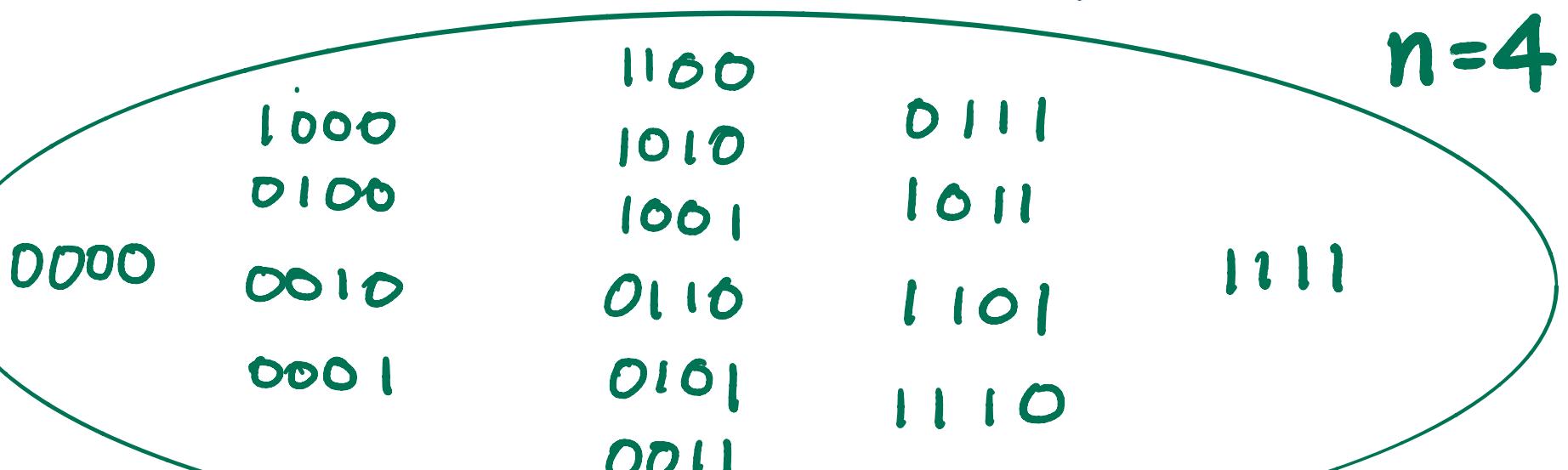
codeword

↑  
code  
space

max # msgs s.t. Pr. error = 0?







max # msgs s.t. Pr. error = 0?

0000

1000

0100

0010

0001

1100

1010

1001

0110

0101

0011

0111

1011

1101

1110

1111

0000

1000

0100

0010

0001

1100

1010

1001

0110

0101

0011

0111

1011

1101

1110

1111

0000

1000

0100

0010

0001

1100

1010

1001

0110

0101

0011

0111

1011

1101

1110

1111

0000

1000

0100

0010

0001

1100

1010

1001

0110

0101

0011

0111

1011

1101

1110

1111

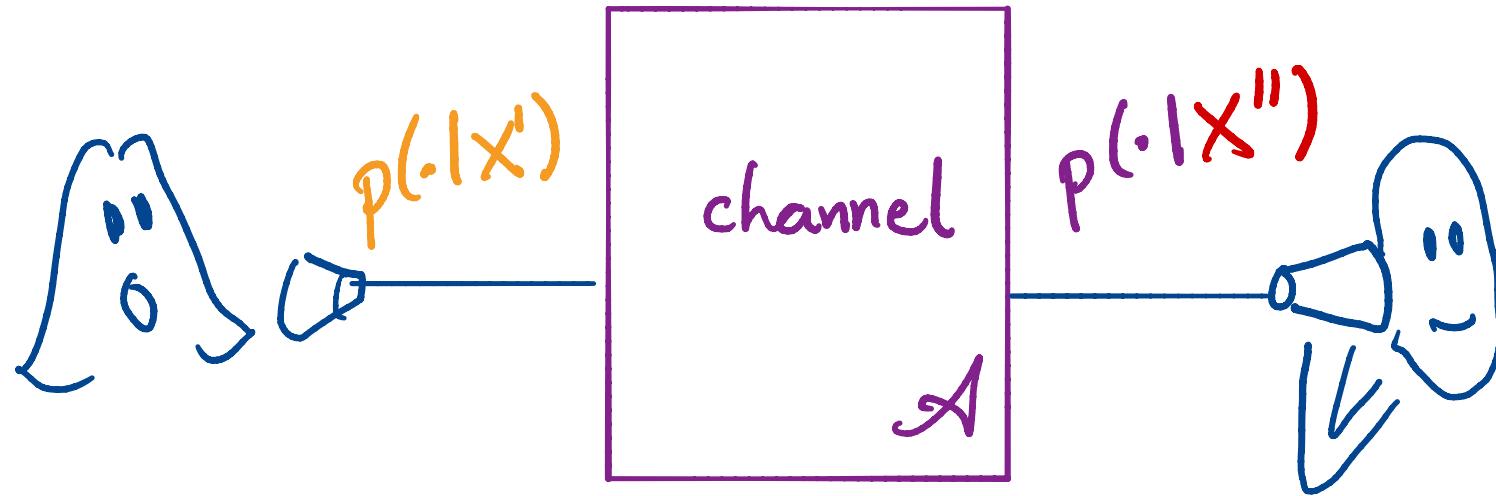
For a noisy channel,  
the max # of msgs  
is defined by the

capacity

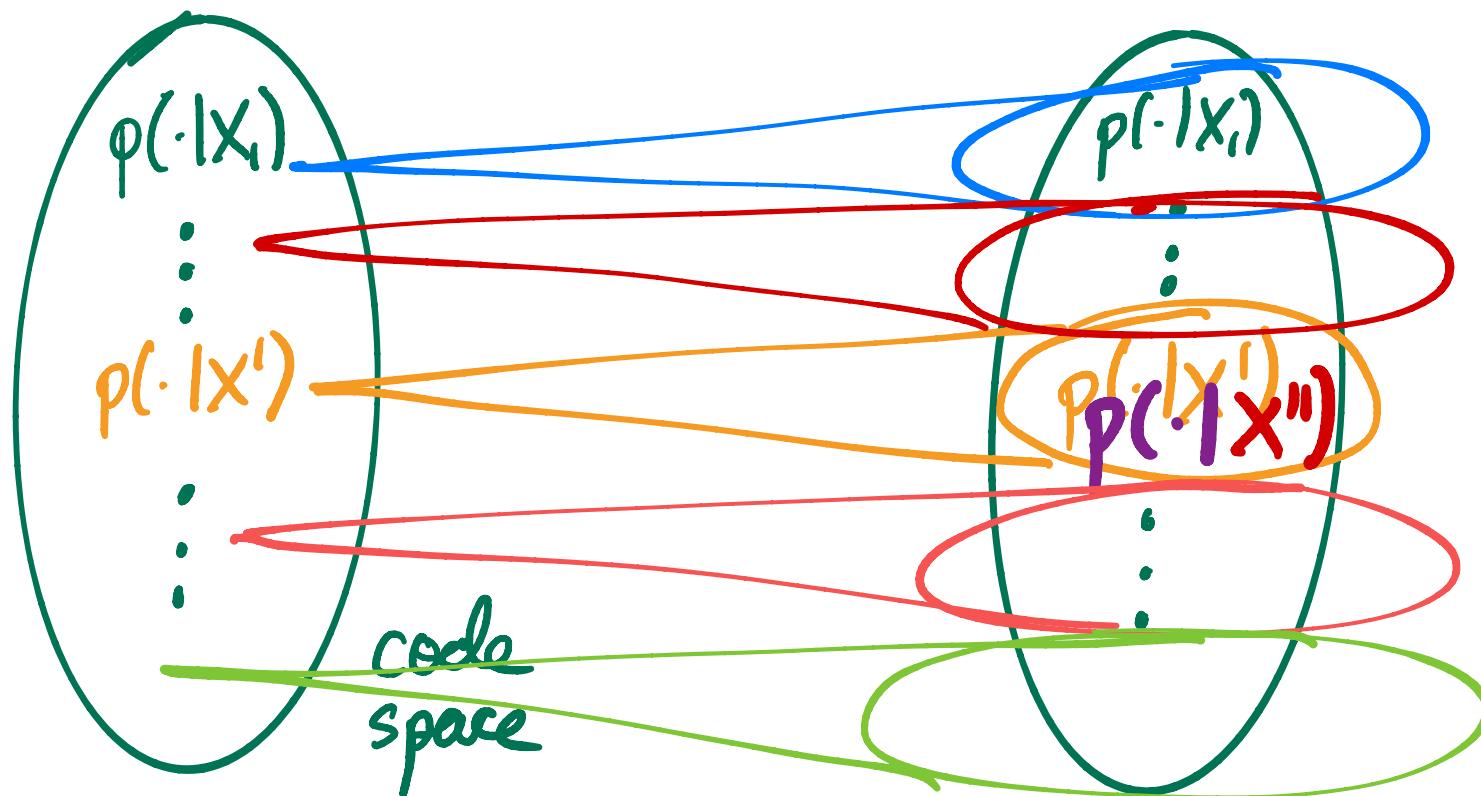
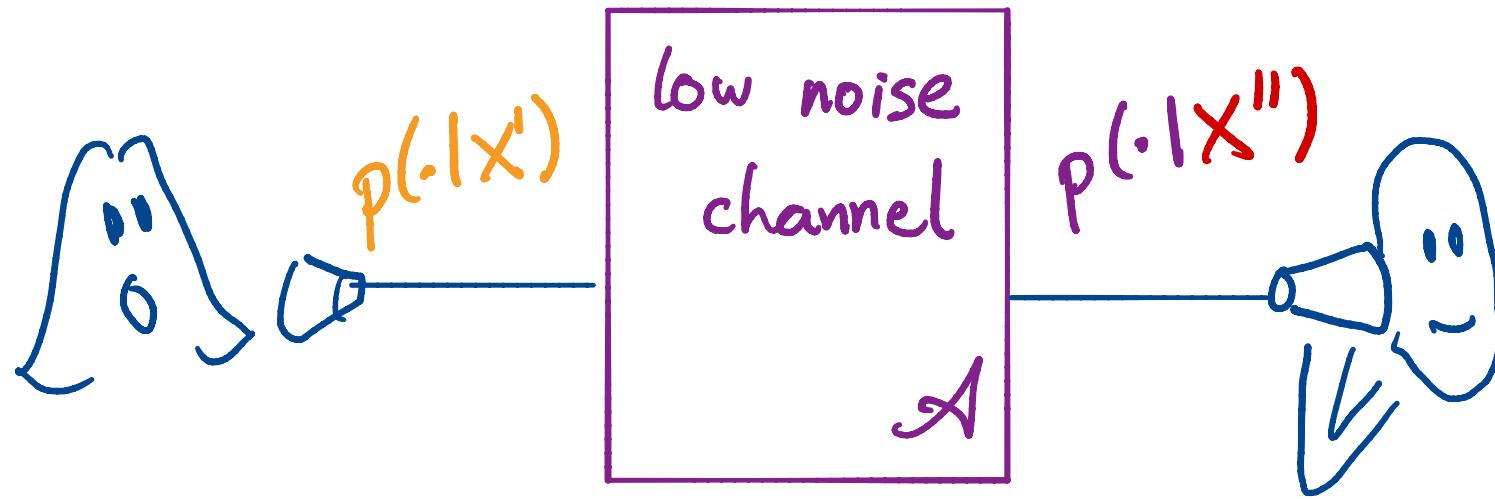
$$\max_p I(S; \hat{S})$$

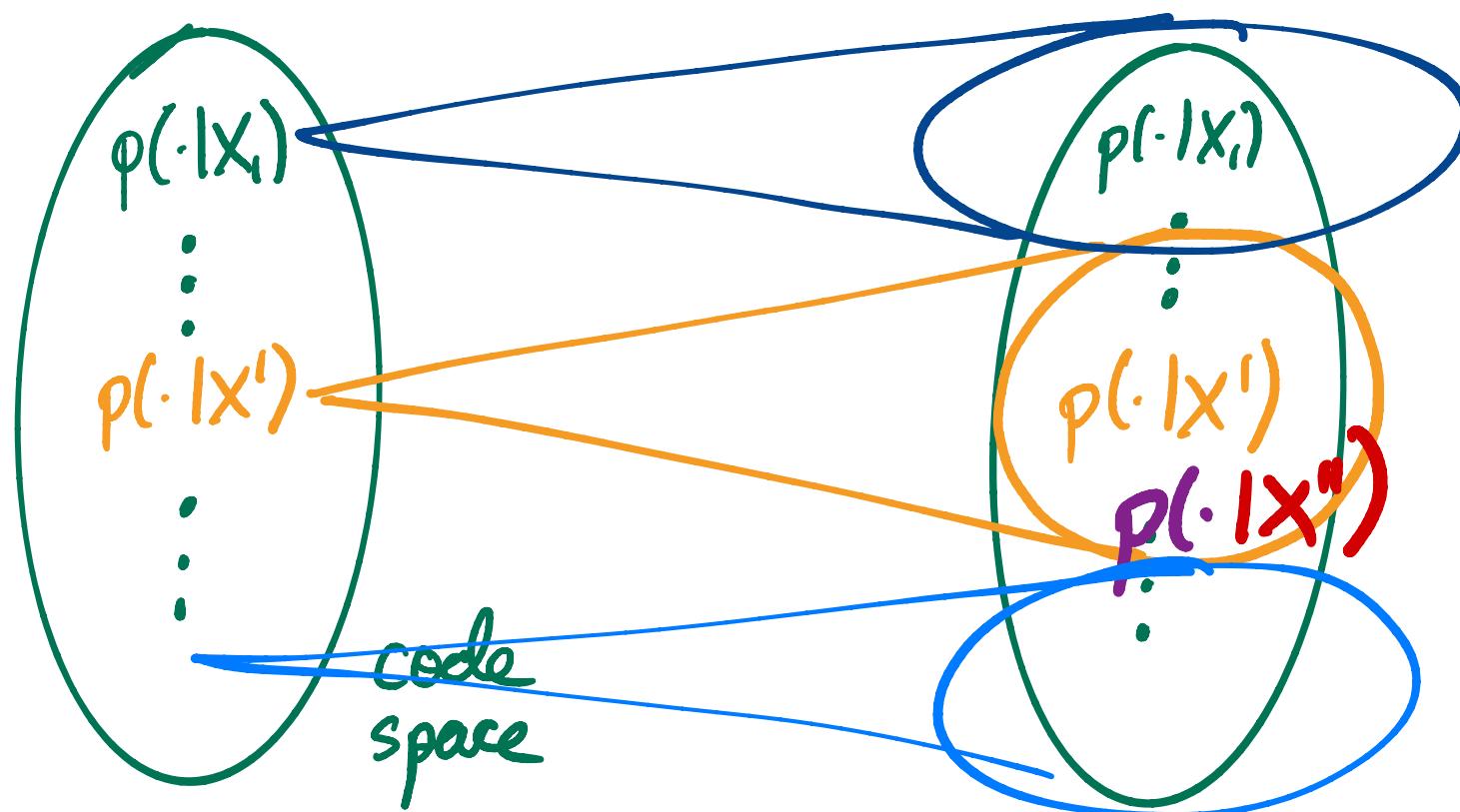
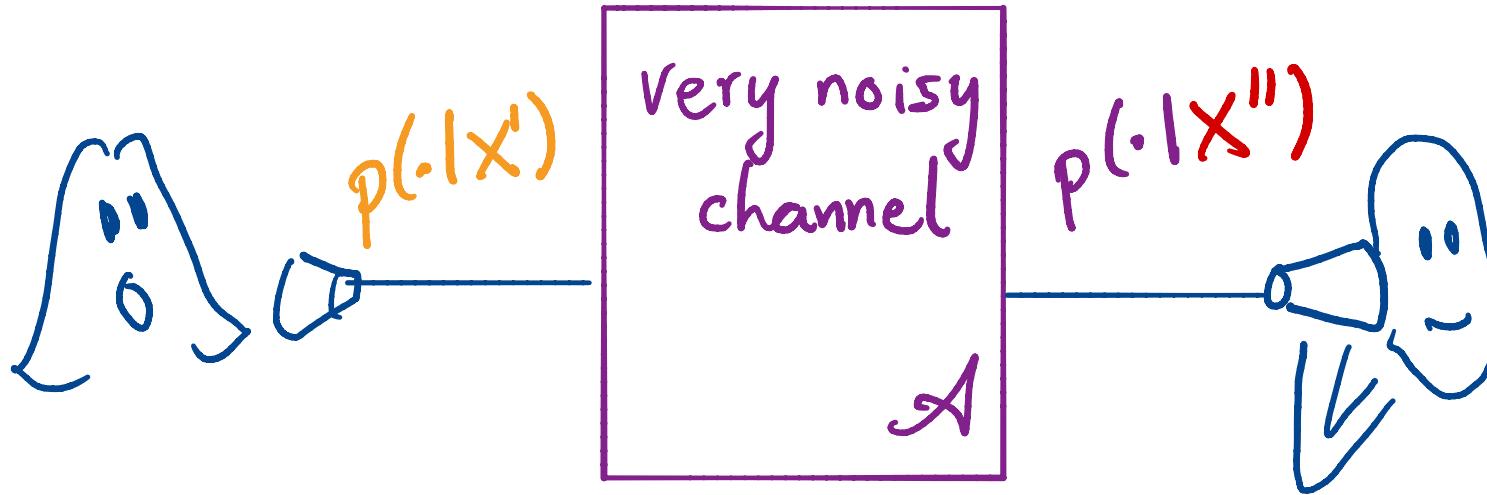
The max # of msgs  
is  $2^{nc}$

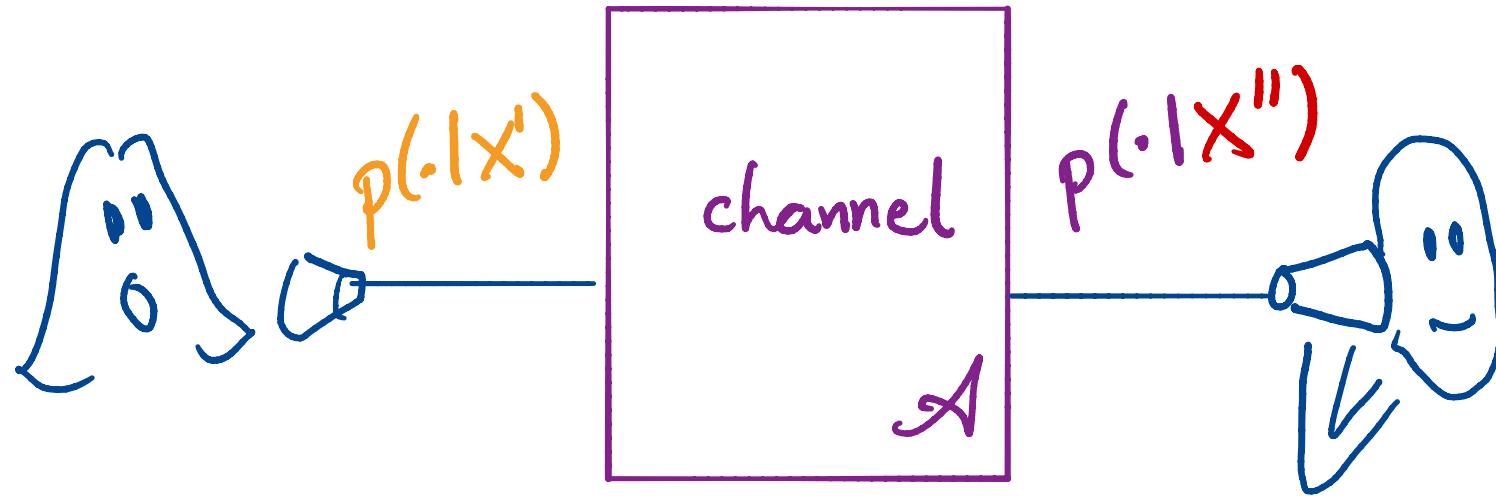
A channel can be  
evaluated by its  
channel capacity



Algorithms are also channels!







An algo is evaluated by its  
"channel's capacity"

## Shannon's coding thm

Aim for channels with high capacity

More capacity  
⇒ more messages

## Posterior agreement

Aim for algorithms with high exp. log. post. agr.

Higher exp. log. PA  
⇒ more messages

# Organization

What is PA?

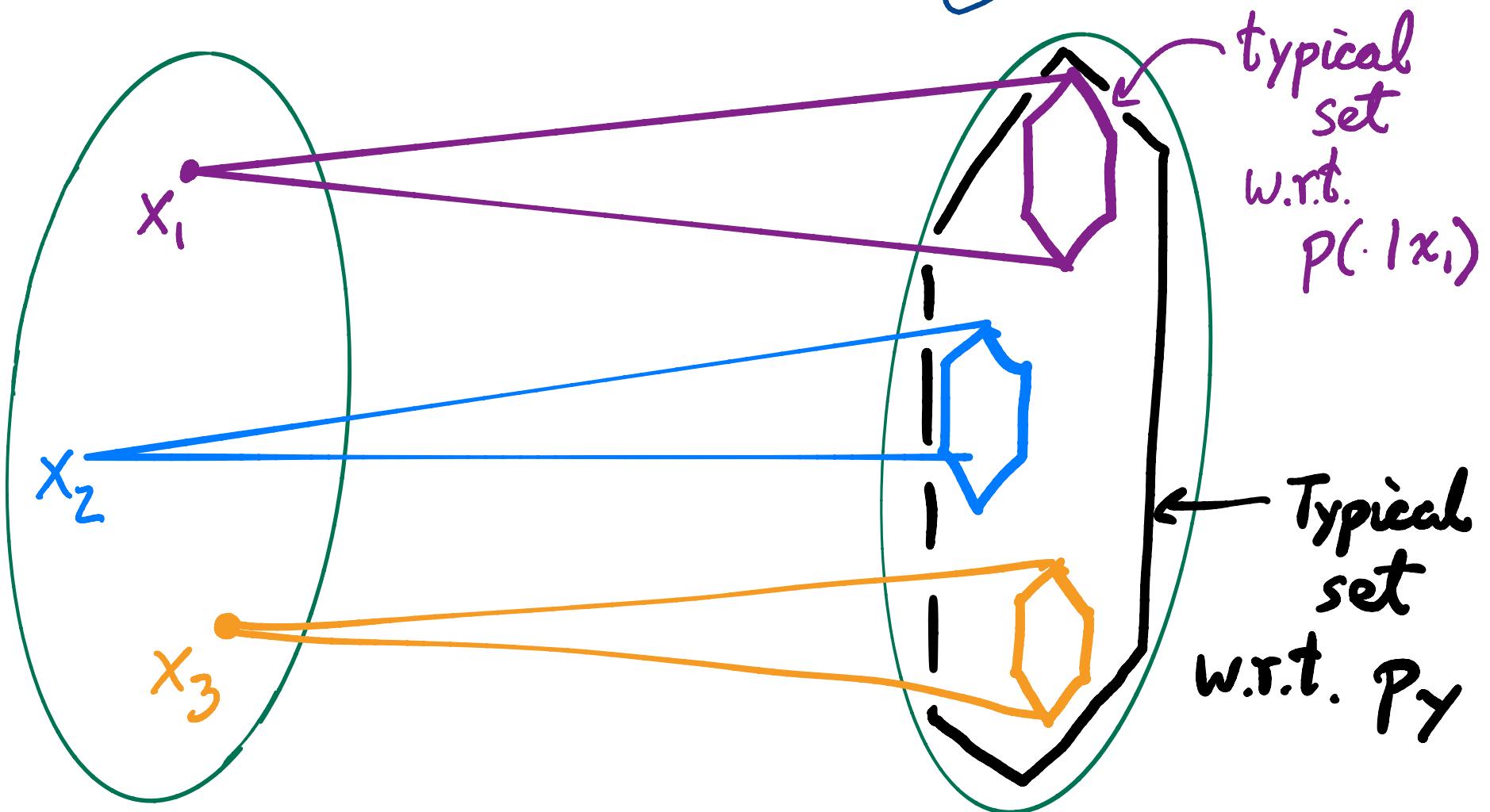
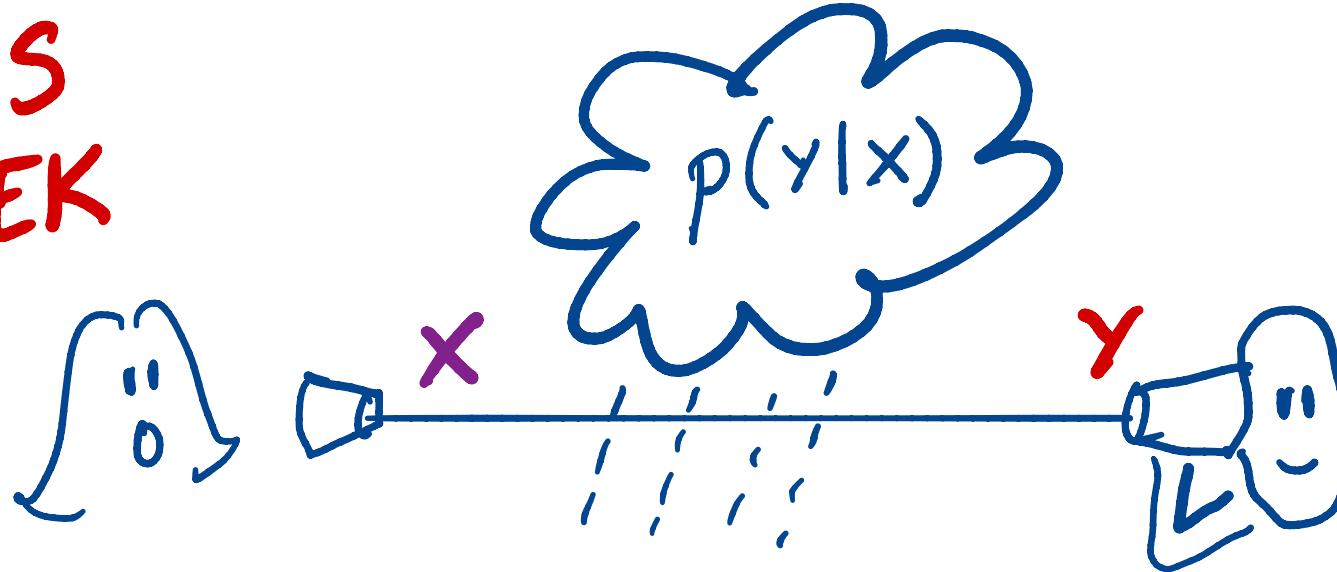
Rationale of PA

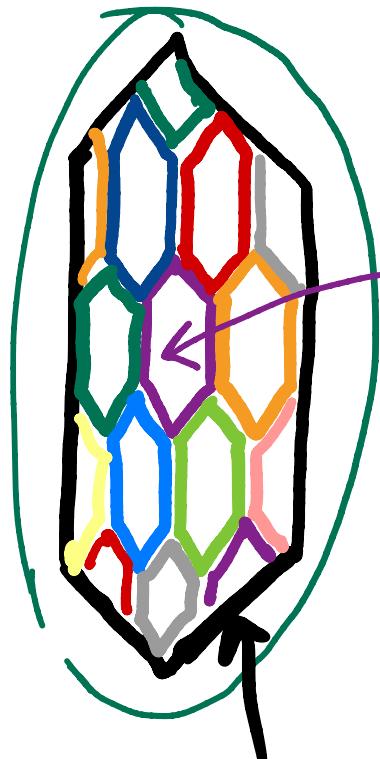
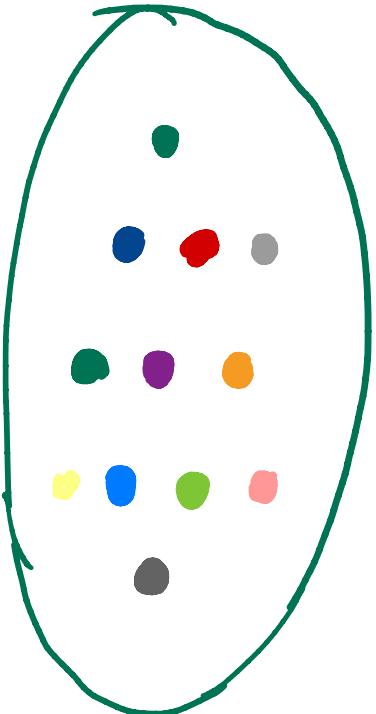
Roadmap

Shannon's coding theorem

Formalization of PA

THIS  
WEEK



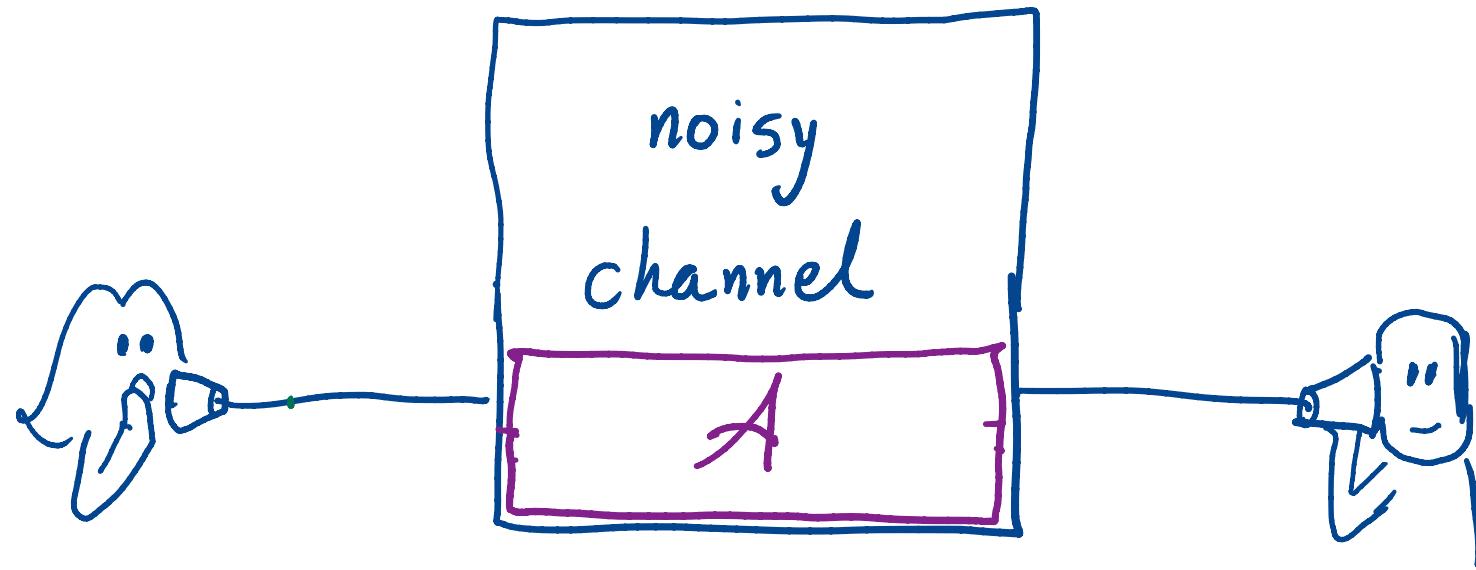


Typical set

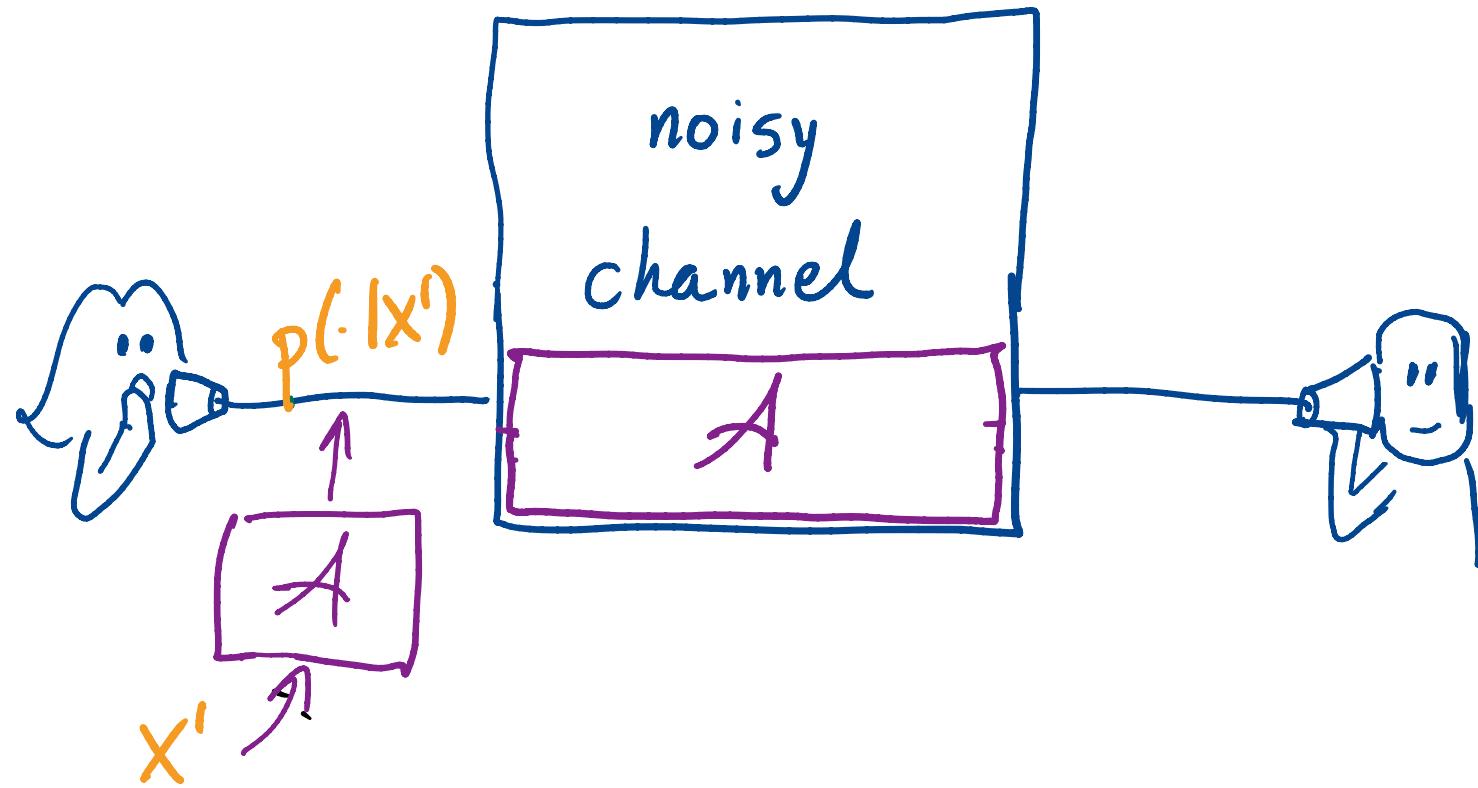
$$\approx 2^{nH(Y|X)}$$

$$\text{Typical set} \approx 2^{nH(Y)}$$

$$\begin{aligned}\text{Optimal \# of msgs} &\leq \frac{2^{nH(Y)}}{2^{nH(Y|X)}} = 2^{nI(X;Y)} \\ &= 2^{nC}\end{aligned}$$

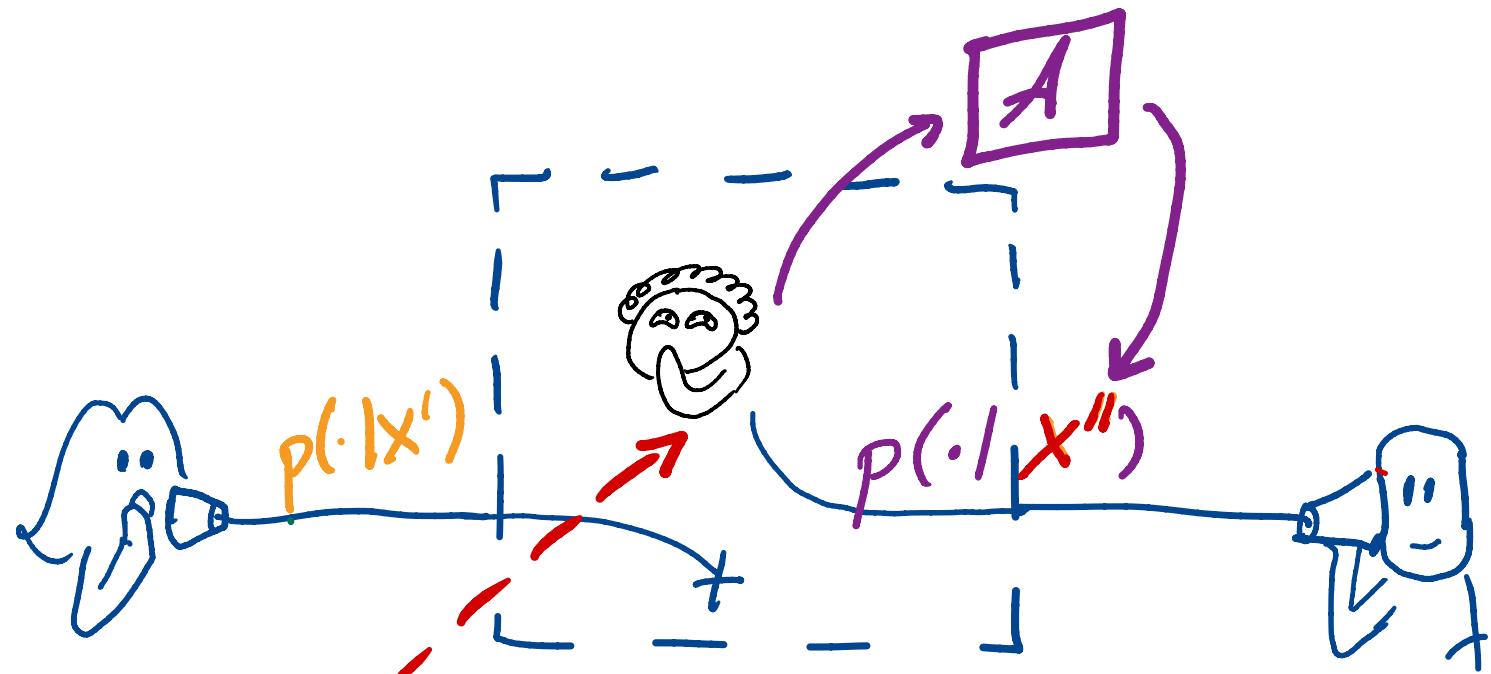


NEXT WEEK

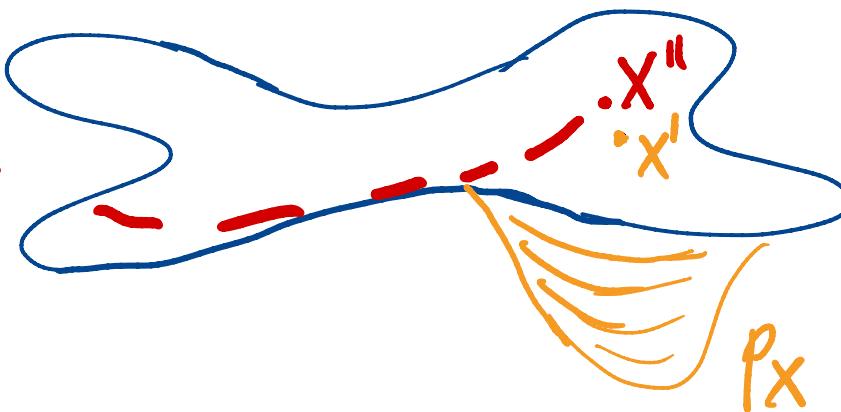


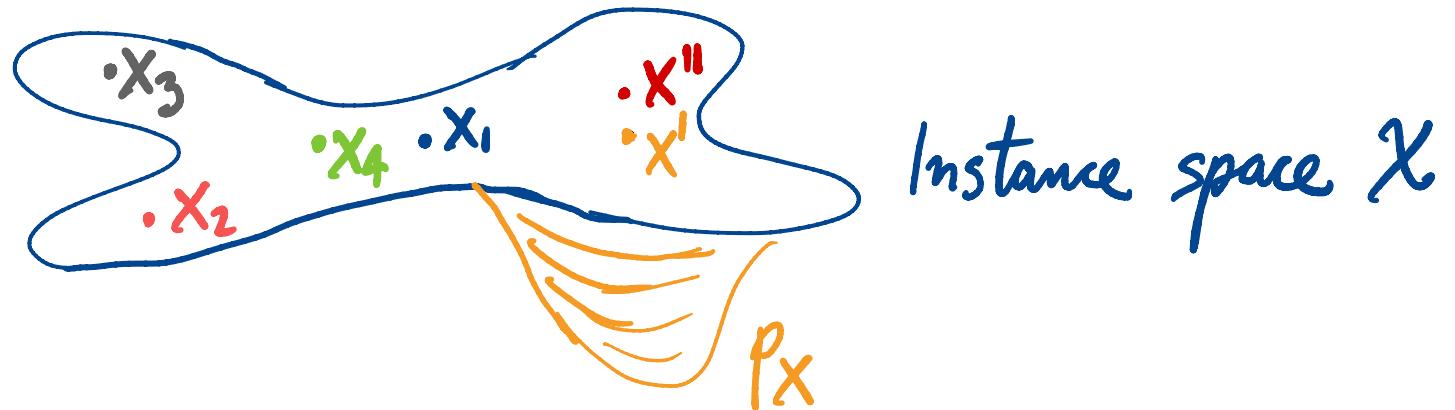
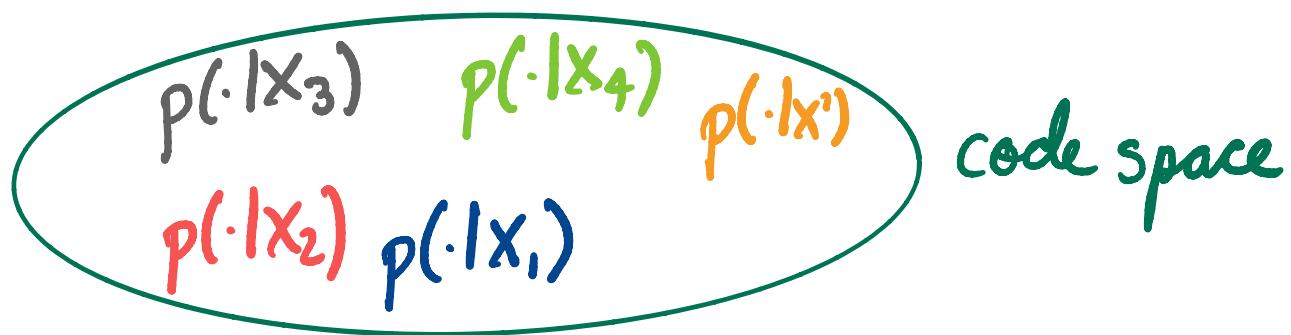
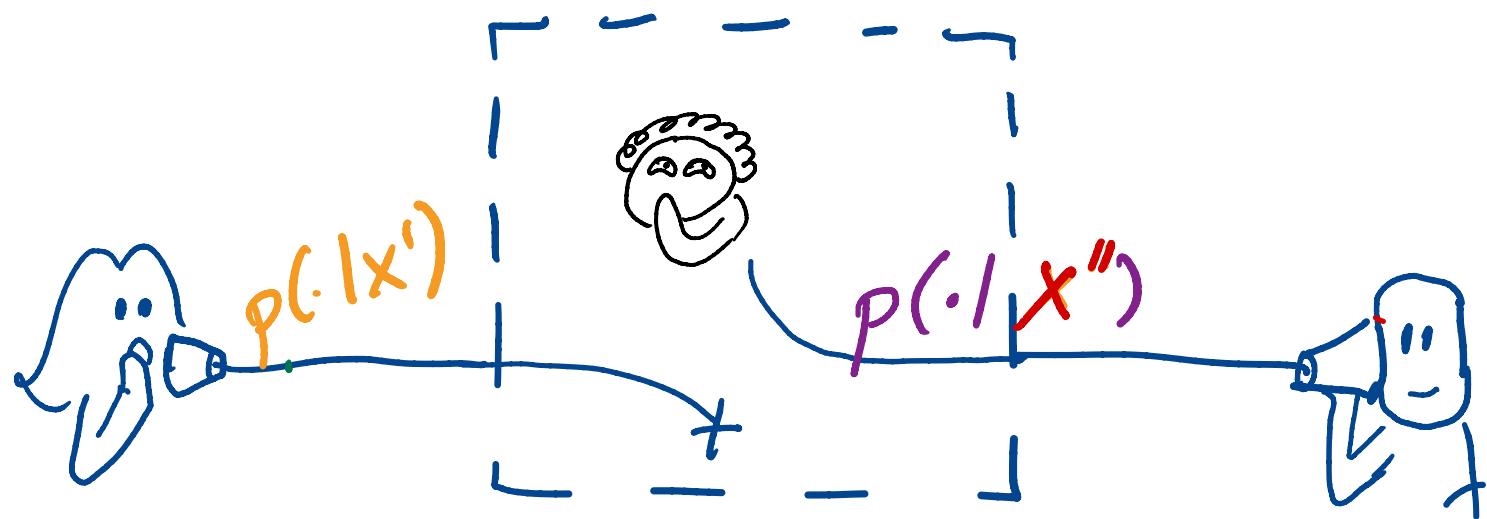
Instance space  $\mathcal{X}$   
 (all instances of "size"  $n$ )

NEXT WEEK  
 $p_{\mathcal{X}} \leftarrow \text{experiment}$

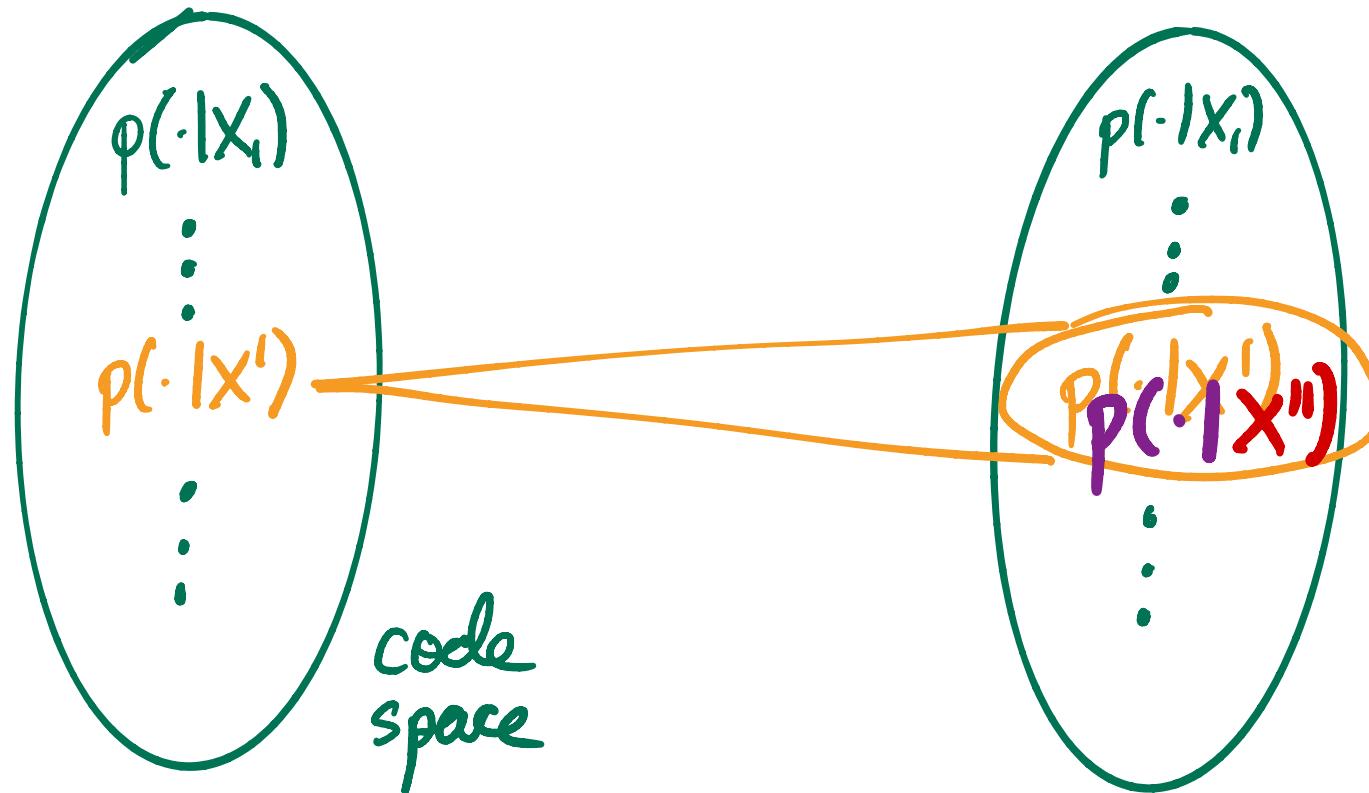
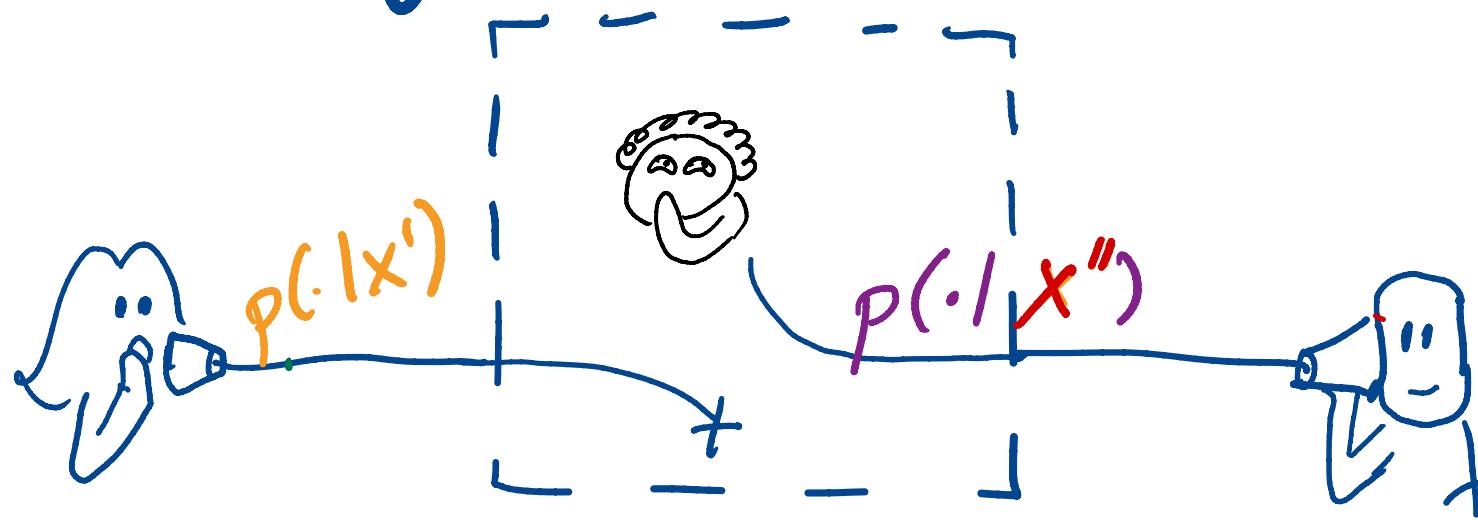


$x''$

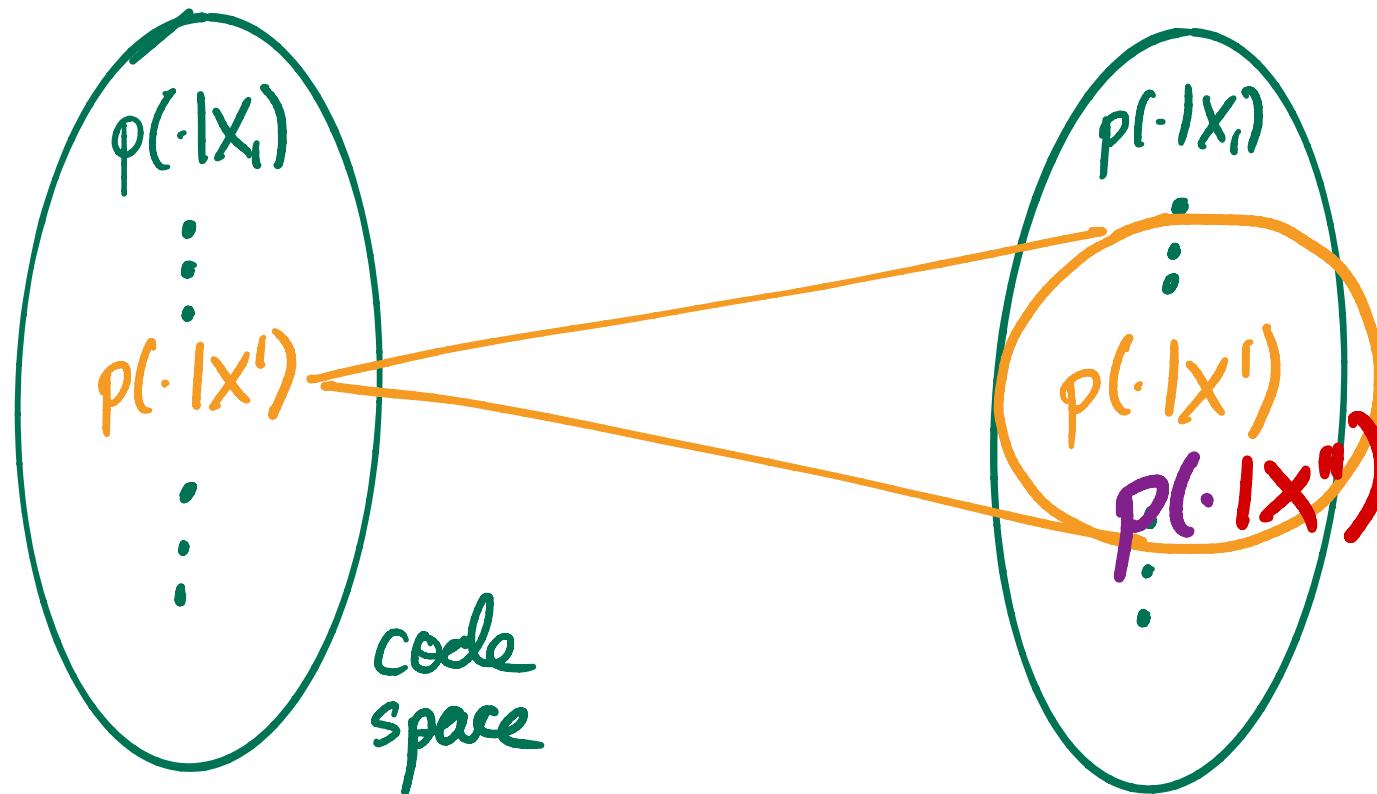
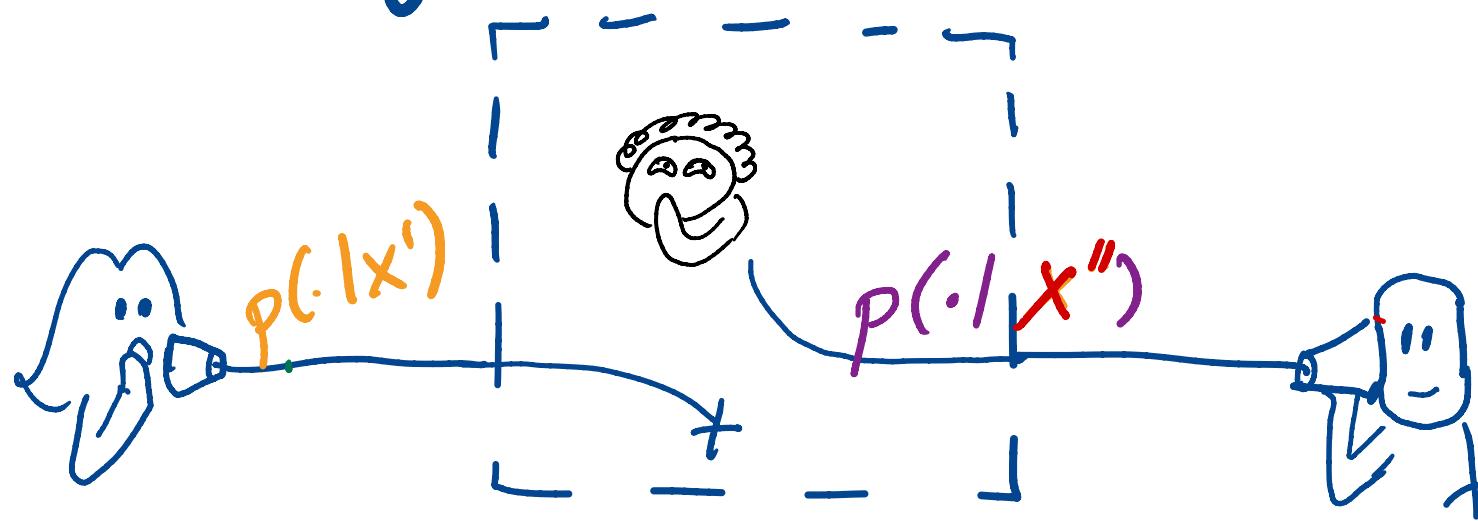


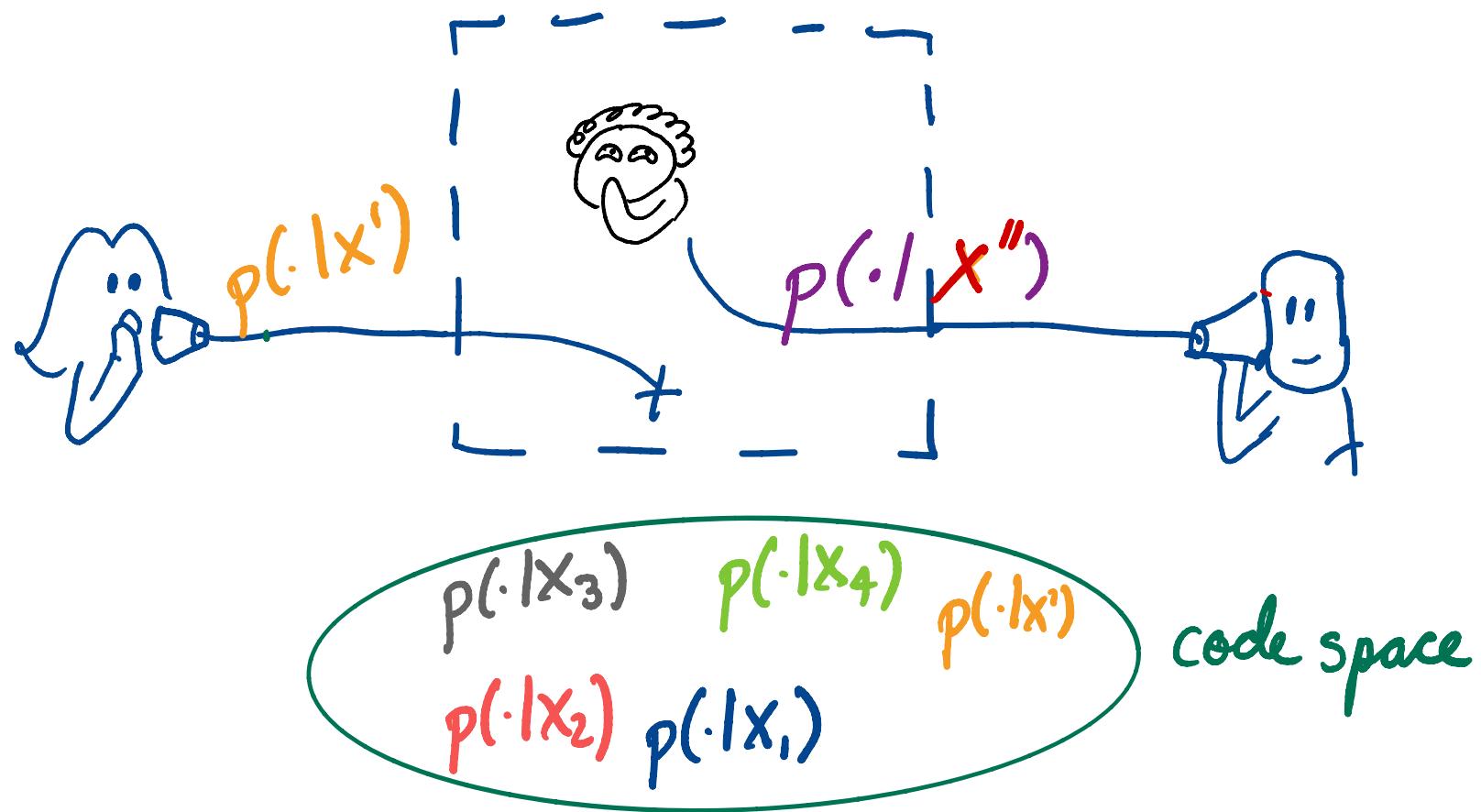


If the algo is robust...



If the algo is not robust...





Channel capacity = exp. log. post. agr.

We prove later that, with high probability,

prob fail  $\leq$

# of msgs  
↓

$$\text{const} \cdot \exp\left(-\log |C| \left(\exp. \log PA - \frac{\log m_n}{\log |C|} - \varepsilon\right)\right)$$

So, if  $\exp. \log PA > \frac{\log m_n}{\log |C|}$  then

prob. fail  $\xrightarrow[n \rightarrow \infty]{\text{in prob.}} 0$

$$\text{const} \cdot \exp\left(-\log |C| \left(\exp. \log PA - \frac{\log m_n}{\log |C|} - \epsilon\right)\right)$$

An algorithm can reduce prob of fail by:

maximizing exp. log. post. agreement.

This allows a large number  $m_n$  of msgs  
that sender can communicate to the receiver.

# Organization

What is PA?

Rationale of PA

Roadmap

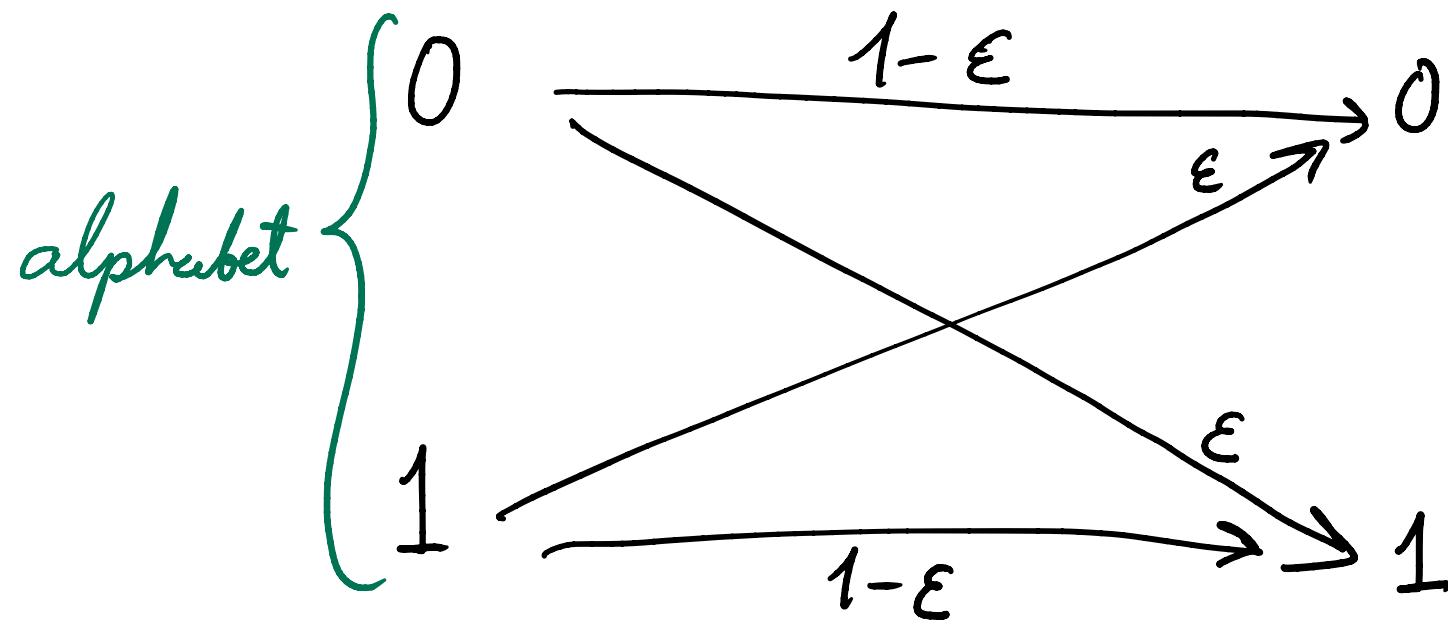
Shannon's coding theorem

Formalization of PA

Shannon's channel

coding theorem

## Channels

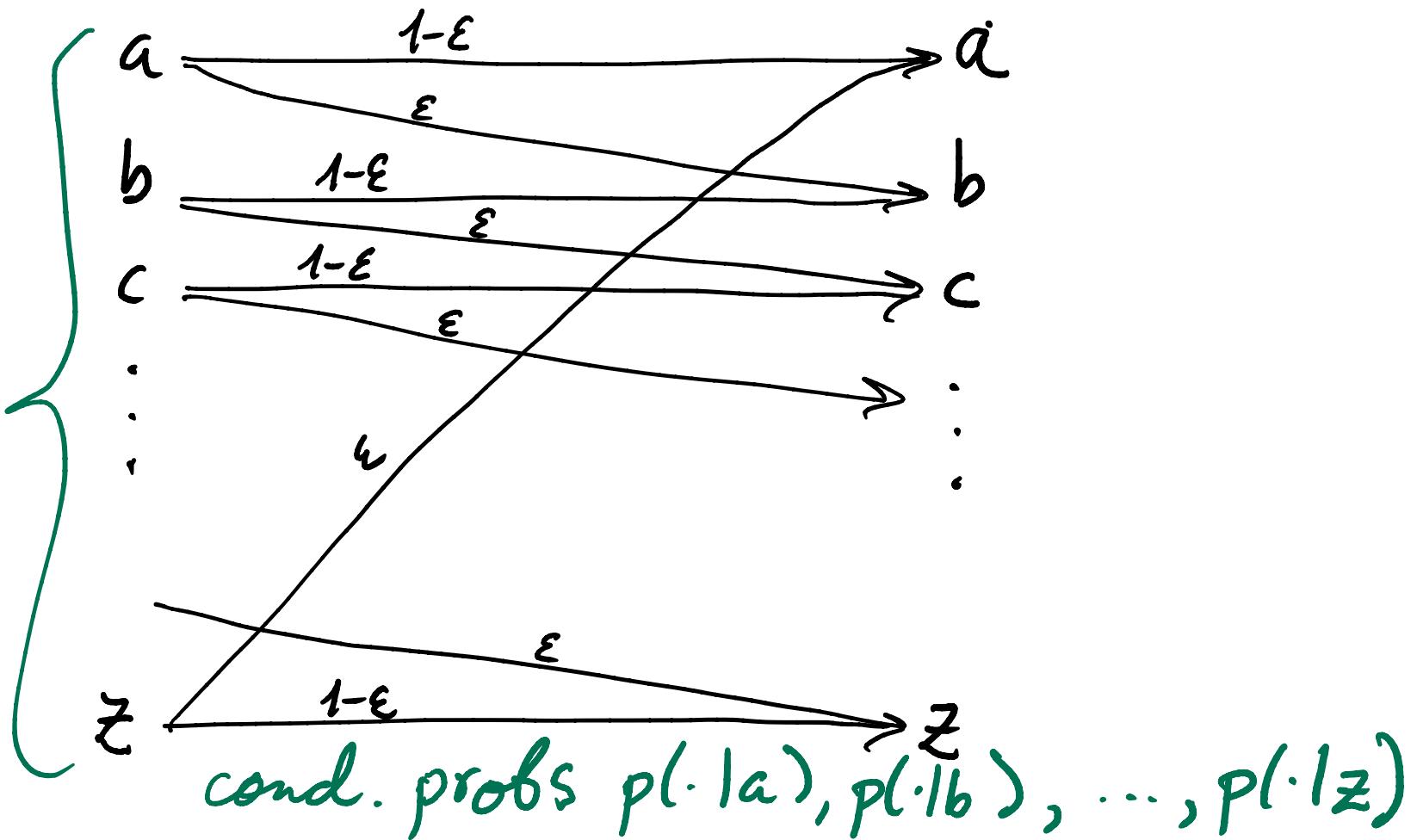


condi. probs  $p(\cdot|0), p(\cdot|1)$

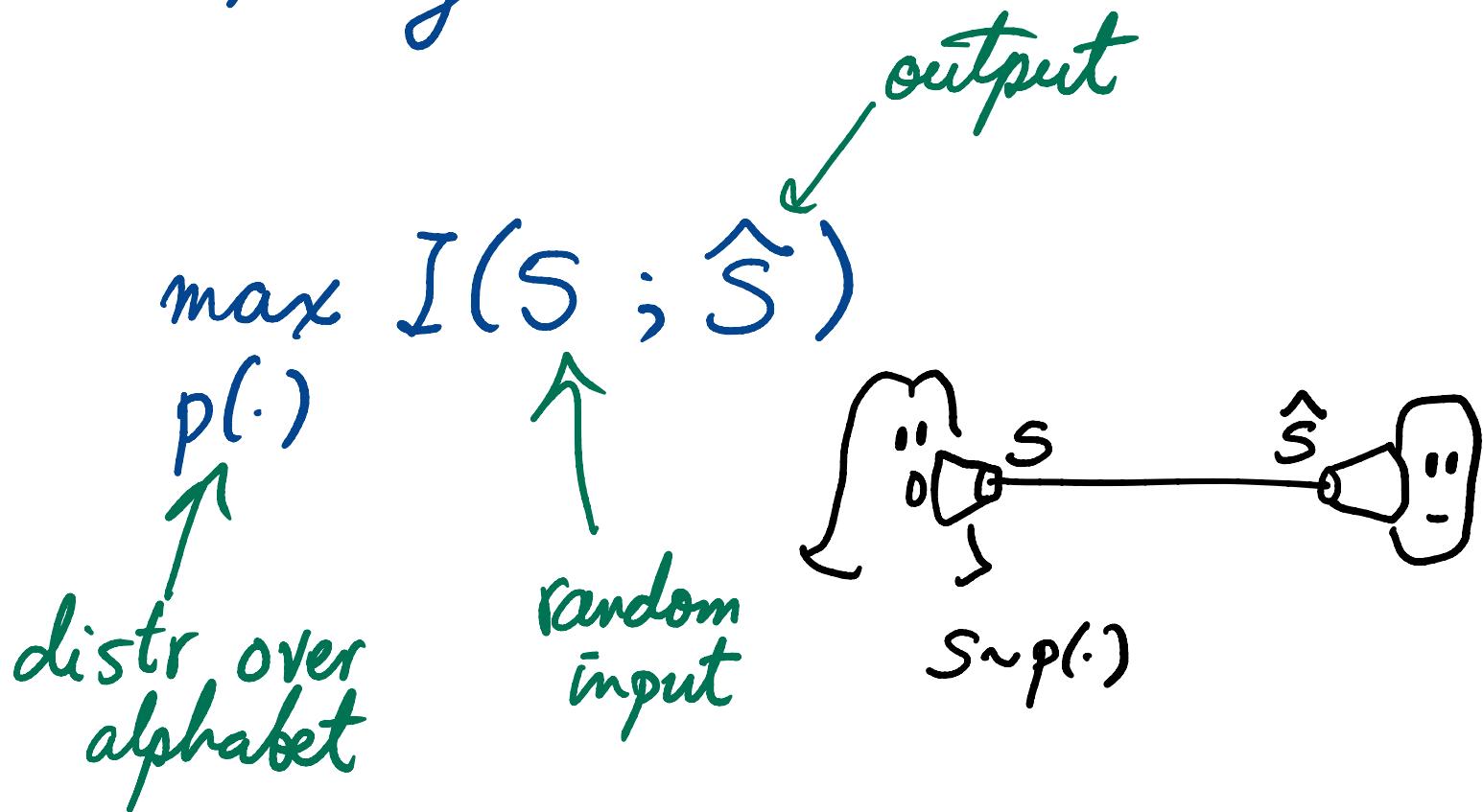
## Channels



alphabet



# Channel capacity



Observe: Typewriter's capacity > binary channel's capacity

## Channel capacity

- \* The max rate in bits / transmission with arbitrarily low prob. of error.
- \* max # of different messages that can be communicated =  $2^{n \text{Capacity}}$

## Quick recap

$$H(X) := - \sum_x p(x) \log p(x)$$

"Info in X"

$$H(X|Y) := - \sum_y p(y) \sum_x p(x|y) \log p(x|y)$$

"residual info left in X  
after learning Y"

$$I(X;Y) := \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

"Info that Y  
reveals about X"

$$H(X) = ?$$

## Quick recap

$$H(X) := -\sum_x p(x) \log p(x)$$

"Info in X"

$$H(X|Y) := -\sum_y p(y) \sum_x p(x|y) \log p(x|y)$$

"residual info left in X  
after learning Y"

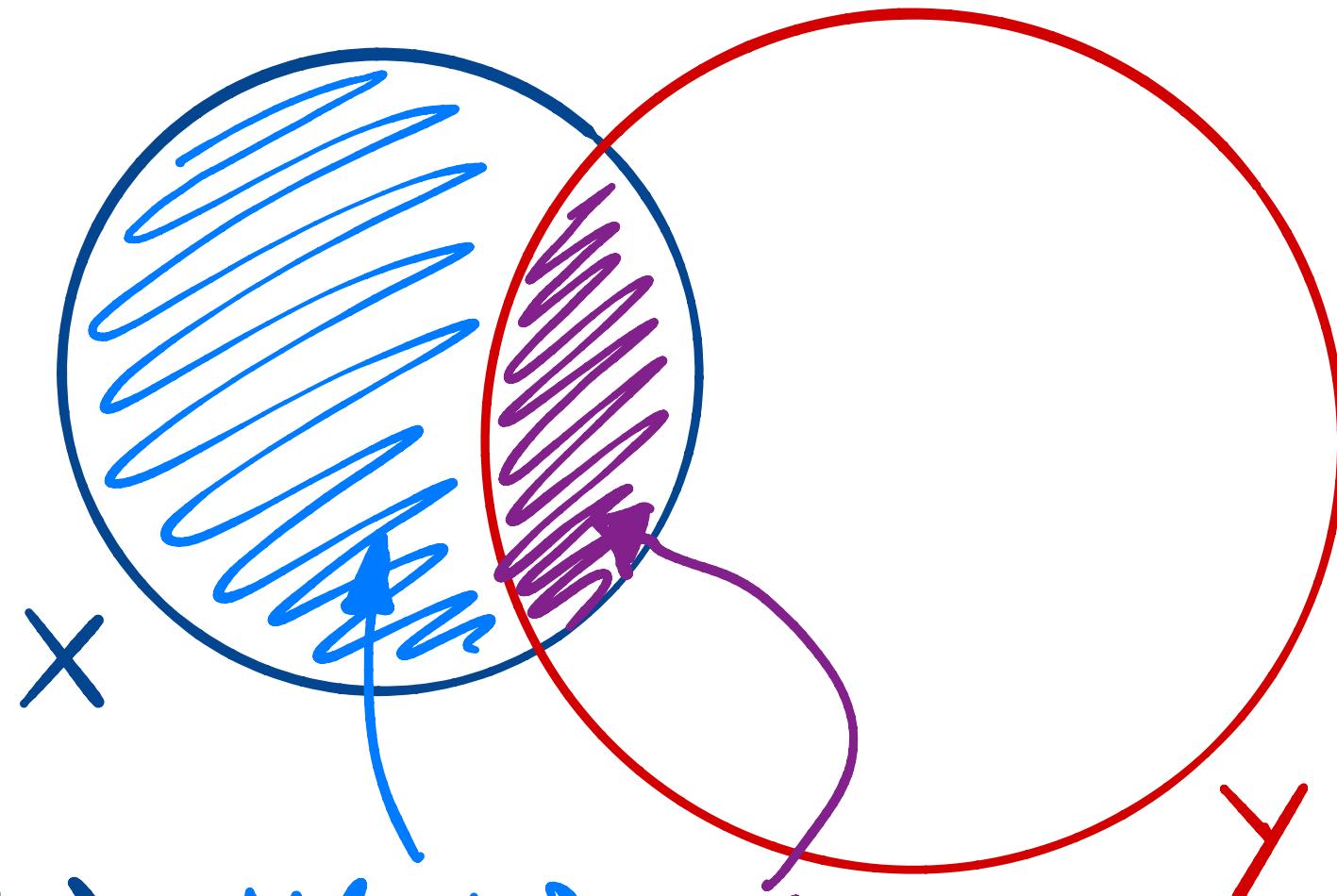
$$I(X;Y) := \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

"Info that Y  
reveals about X"

$$H(X) = I(X;Y) + H(X|Y)$$

$$I(X;Y) = H(X) - H(X|Y)$$

$$I(X;Y) = H(X) - H(X|Y)$$



$$H(X) = H(X|Y) + I(X;Y)$$

Proof:

$$I(X;Y) = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

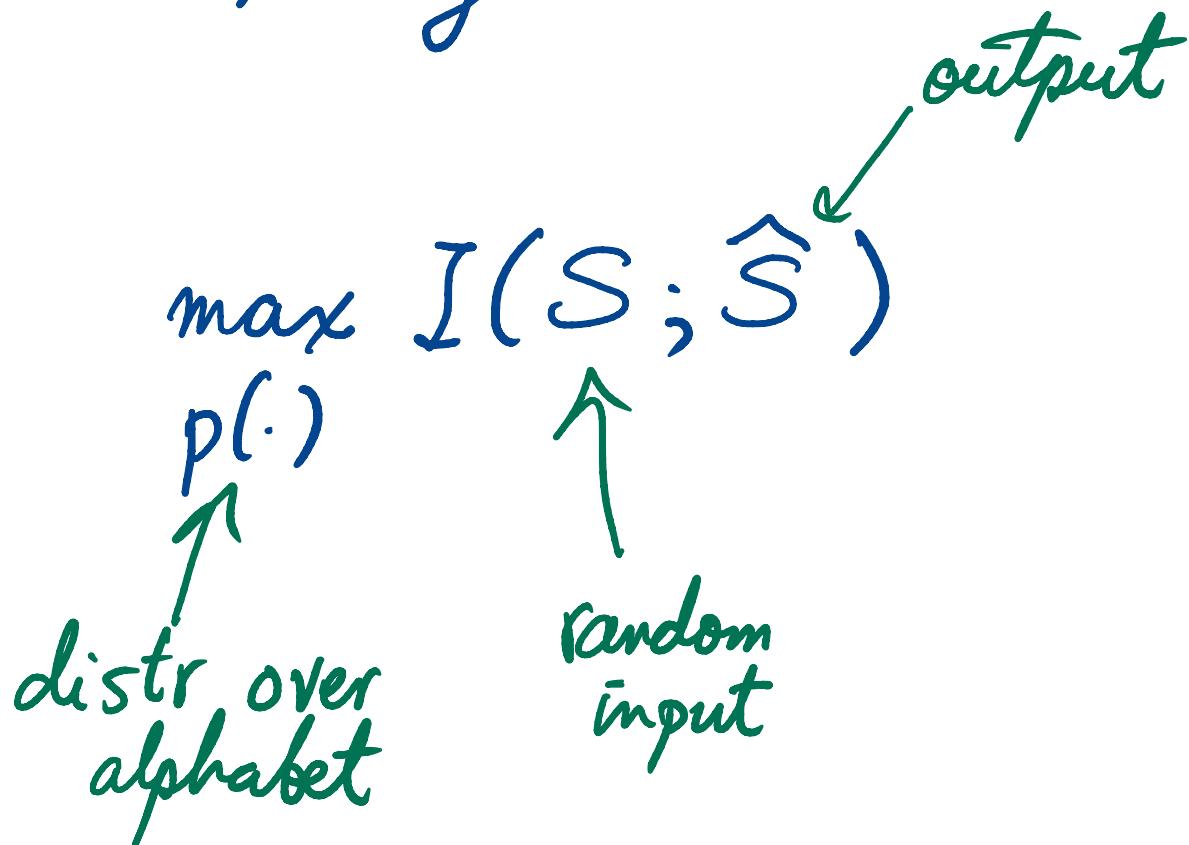
$$= \sum_{x,y} p(x,y) \log \frac{p(x|y)p(y)}{\cancel{p(x)p(y)}}$$

$$= \sum_{x,y} \cancel{p(x,y)} \log p(x|y) - \sum_{x,y} \cancel{p(x,y)} \log p(x)$$

$$= \sum_y p(y) \underbrace{\sum_x p(x|y) \log p(x|y)}_{H(X|Y)} - \sum_x p(x) \underbrace{\sum_y p(y|x) \log p(y|x)}_{H(Y|X)}$$

$$= -H(X|Y) + H(X).$$

## Channel capacity



Observe: Typewriter's capacity > binary channel's capacity

Example:

Zero-noise bin channel

$$\max_P I(S; \hat{S}) = ?$$

Example:

Zero-noise bin channel

$$\max_P I(S; \hat{S}) = \max_P H(S) - H(S|\hat{S})$$

$\nearrow^0$   
 $\searrow^1$

$$= \max_P H(S) = 1$$

1 bit per transmission

Zero-noise typewriter

$$\max_P I(S; \hat{S}) = ?$$

Zero-noise typewriter

$$\begin{aligned}\max_P I(S; \hat{S}) &= \max_P H(S) - H(S|\hat{S}) \\ &= \max_P H(S) = \log 26 \approx 4.7\end{aligned}$$

4.7 bits per transmission

Noisy typewriter with  $\epsilon = 0.5$

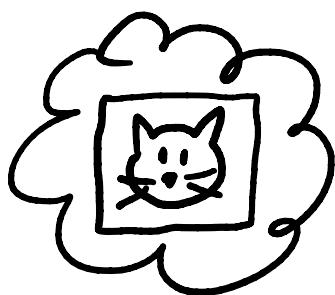
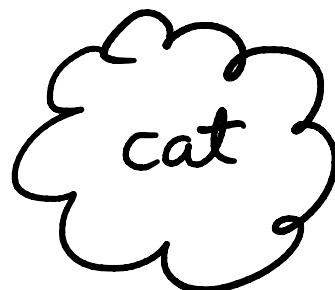
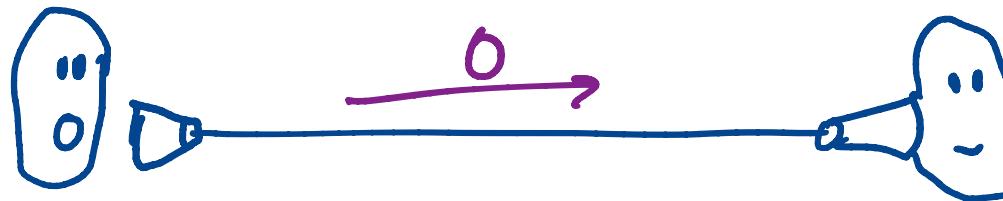
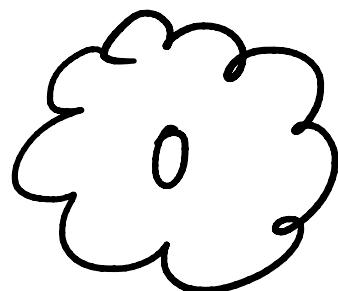
$$\max_P I(S; \hat{S}) = ?$$

Noisy typewriter with  $\epsilon = 0.5$

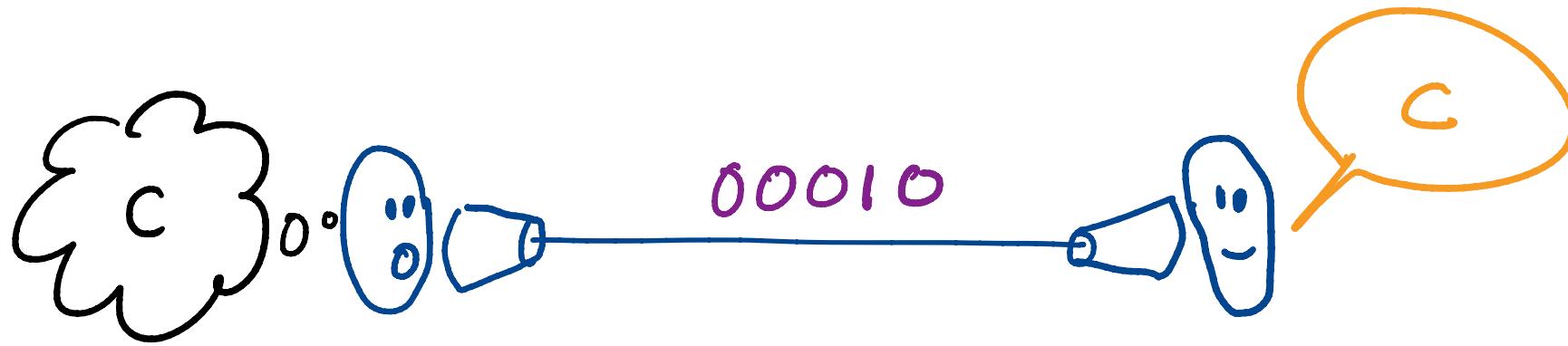
$$\begin{aligned}\max_P I(S; \hat{S}) &= \dots = \log(26/2) \\ &= \log 26 - 1 \\ &\approx 4.7 - 1 \\ &= 3.7\end{aligned}$$

3.7 bits per transmission

Codes:



# Naïve alphabet code



a 00000 ← codewords

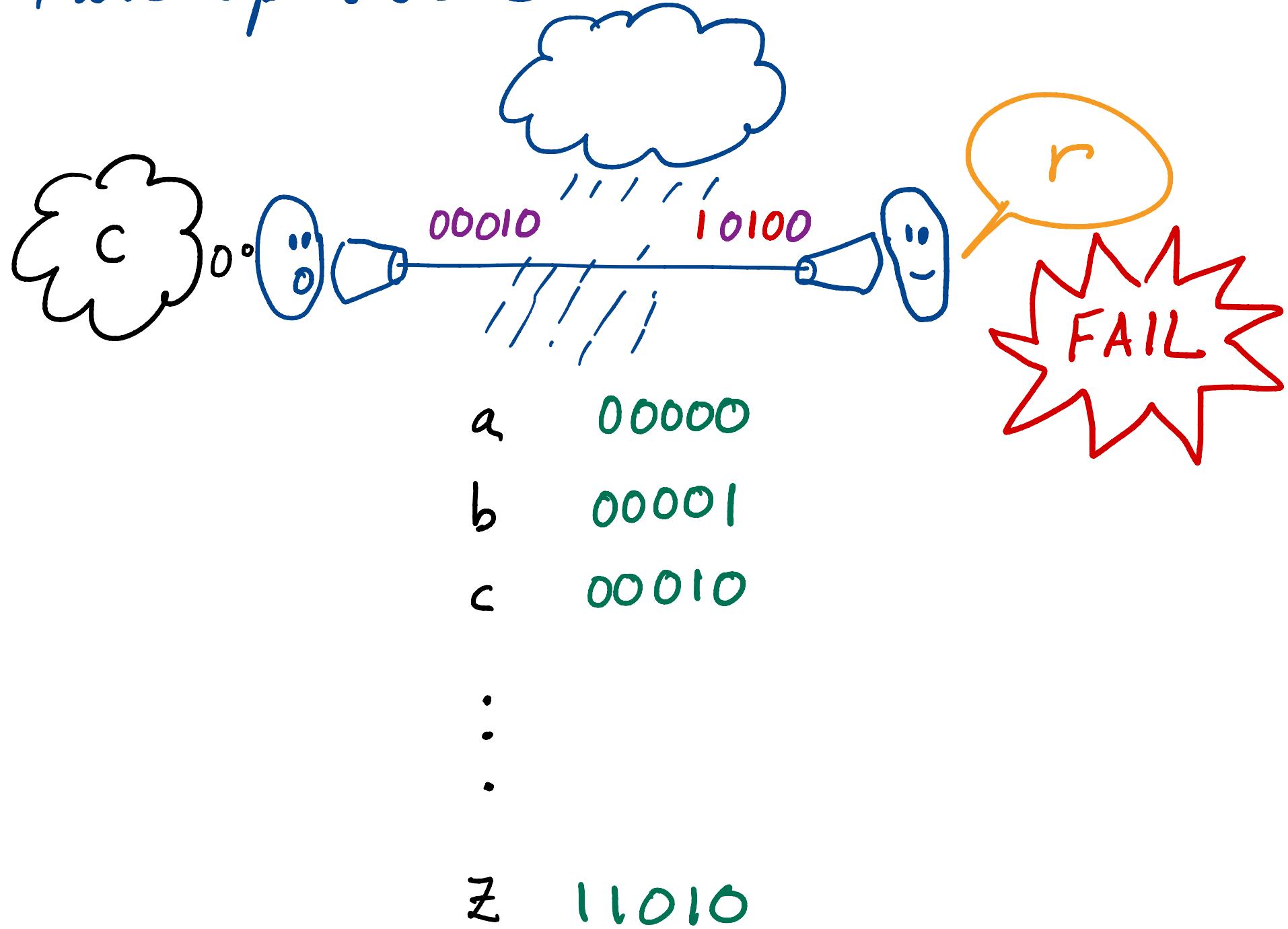
b 00001

c 00010

⋮

z 11001

# Naïve alphabet code



Info theory has come up with  
robust codes...

... but at the cost of  
less bits per transmission

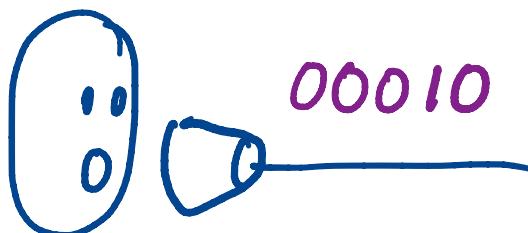
5-redundant code



a 00000

b 11111

Naïve



codeword  
length

$n$

# messages

$2^n$  "

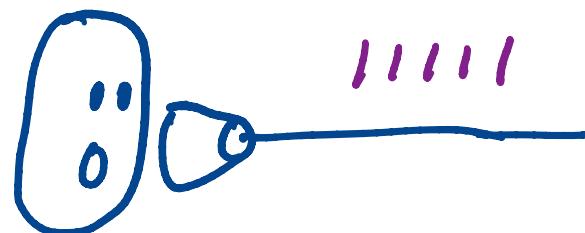
prob. fail  
 $n \rightarrow \infty$

1 "

bits/trans  
 $n \rightarrow \infty$

1 "

$n$ -redundant



$n$

$2$  "

0 "

0 "

Goal: Max. rate as prob. fail  $\rightarrow 0$  as  $n \rightarrow \infty$

Shannon's coding thm

If # of msgs =  $2^{nr}$  (with  $r < C$ )

$\Rightarrow$  there is a code

$$\text{prob. error} \xrightarrow[n \rightarrow \infty]{} 0$$

(The converse is true)

This means  $\Rightarrow$   $\max \# \text{msgs} = 2^{nC}$

$\Rightarrow$  capacity as a channel measure

# Shannon's random code

$$m_1 \quad S_1^n := S_{11}^n S_{12}^n \dots S_{1n}^n$$

$$m_2 \quad S_2^n := S_{21}^n S_{22}^n \dots S_{2n}^n$$

.

.

$$m_{2^{nr}} \quad S_{2^{nr}}^n := \dots$$

# Shannon's random code

$$m_1 \quad S_1^n := S_{11}^n S_{12}^n \dots S_{1n}^n$$

$$m_2 \quad S_2^n := S_{21}^n S_{22}^n \dots S_{2n}^n$$

.

:

$$m_{2^{nr}} \quad S_{2^{nr}}^n := \dots$$

$$S_{i,j}^n \xleftarrow{\$} p_S^* = \arg \max_P I(S; \hat{S})$$

Every symbol of  
every codeword is  
chosen indep. at random  
from  $p_S^* = \arg \max_P I(S; \hat{S})$

# Shannon's random code

$$m_1 \quad S_1^n := S_{11}^n S_{12}^n \dots S_{1n}^n$$

$$m_2 \quad S_2^n := S_{21}^n S_{22}^n \dots S_{2n}^n$$

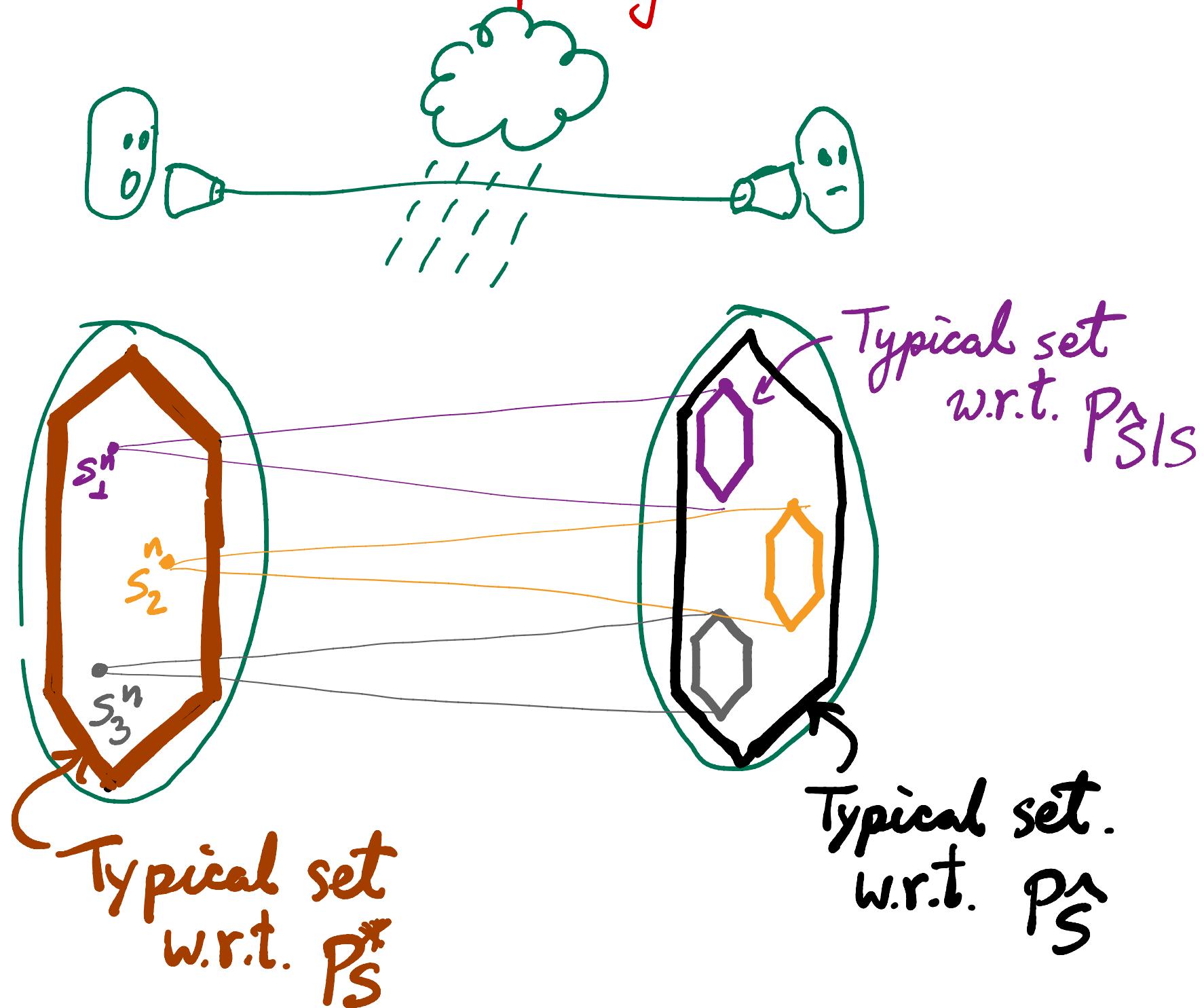
$$\vdots$$
$$m_{2^{nr}} \quad S_{2^{nr}}^n := \dots$$

$$S_{i,j}^n \xleftarrow{\text{symbol}} P_S^* = \arg \max_P I(S; \hat{S})$$

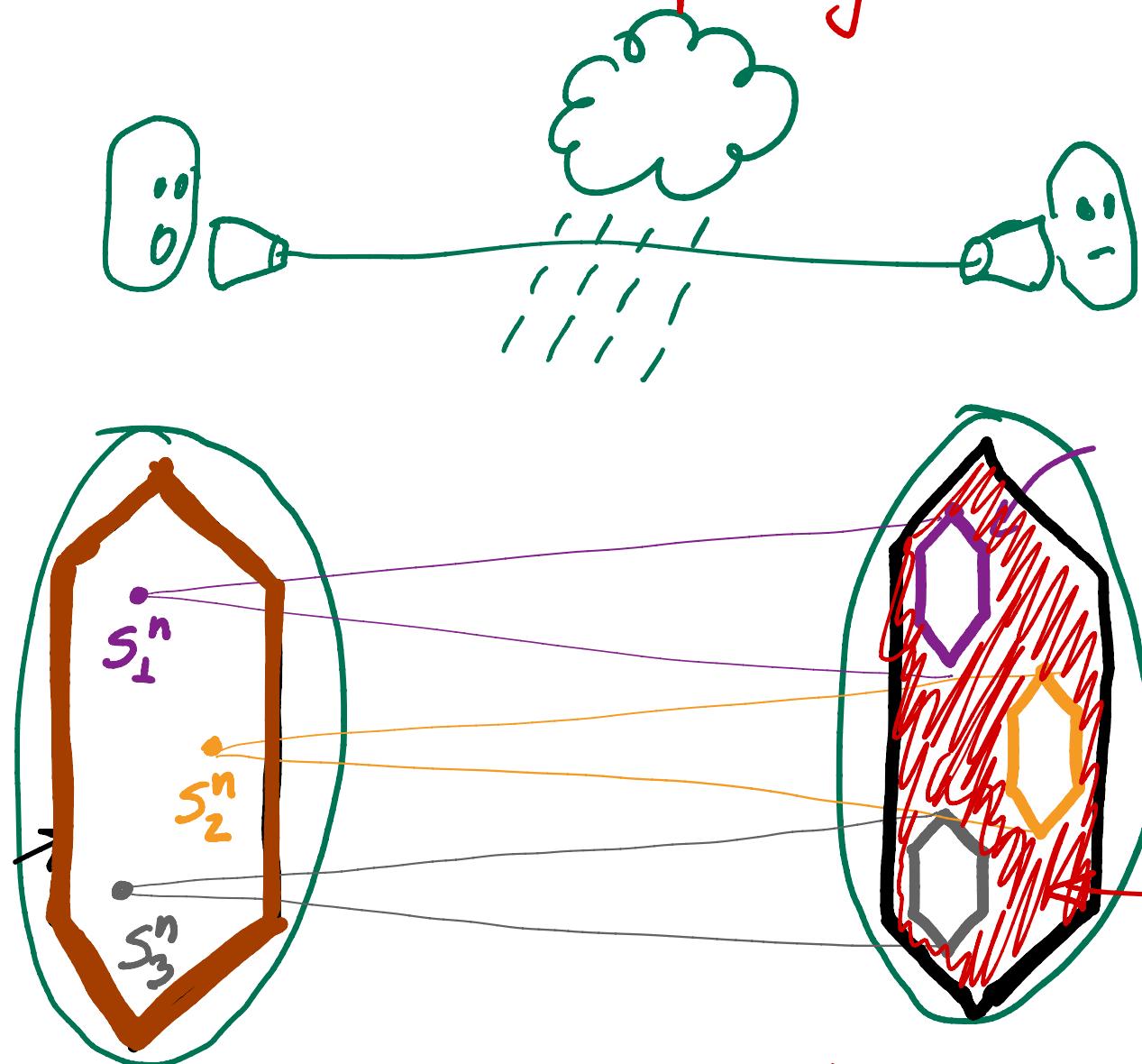
Every symbol of  
every codeword is  
chosen indep. at random  
from  $P_S^* = \arg \max_P I(S; \hat{S})$

$$S_i^n \xleftarrow{\text{codeword}} P_{S^n}^* := \prod_{j \leq n} P_S^*$$

What's the max # of msgs we can communicate?

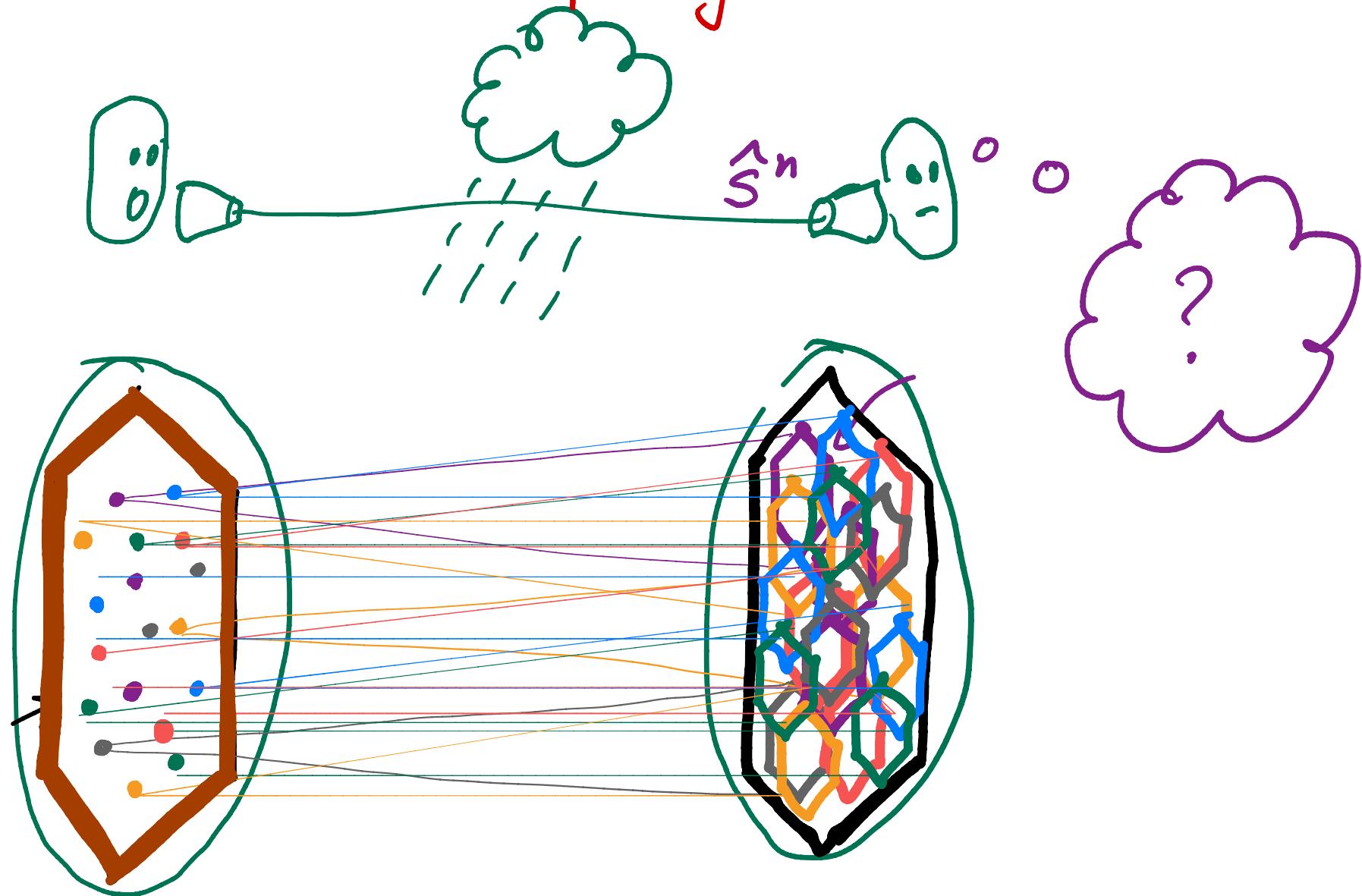


What's the max # of msgs we can communicate?



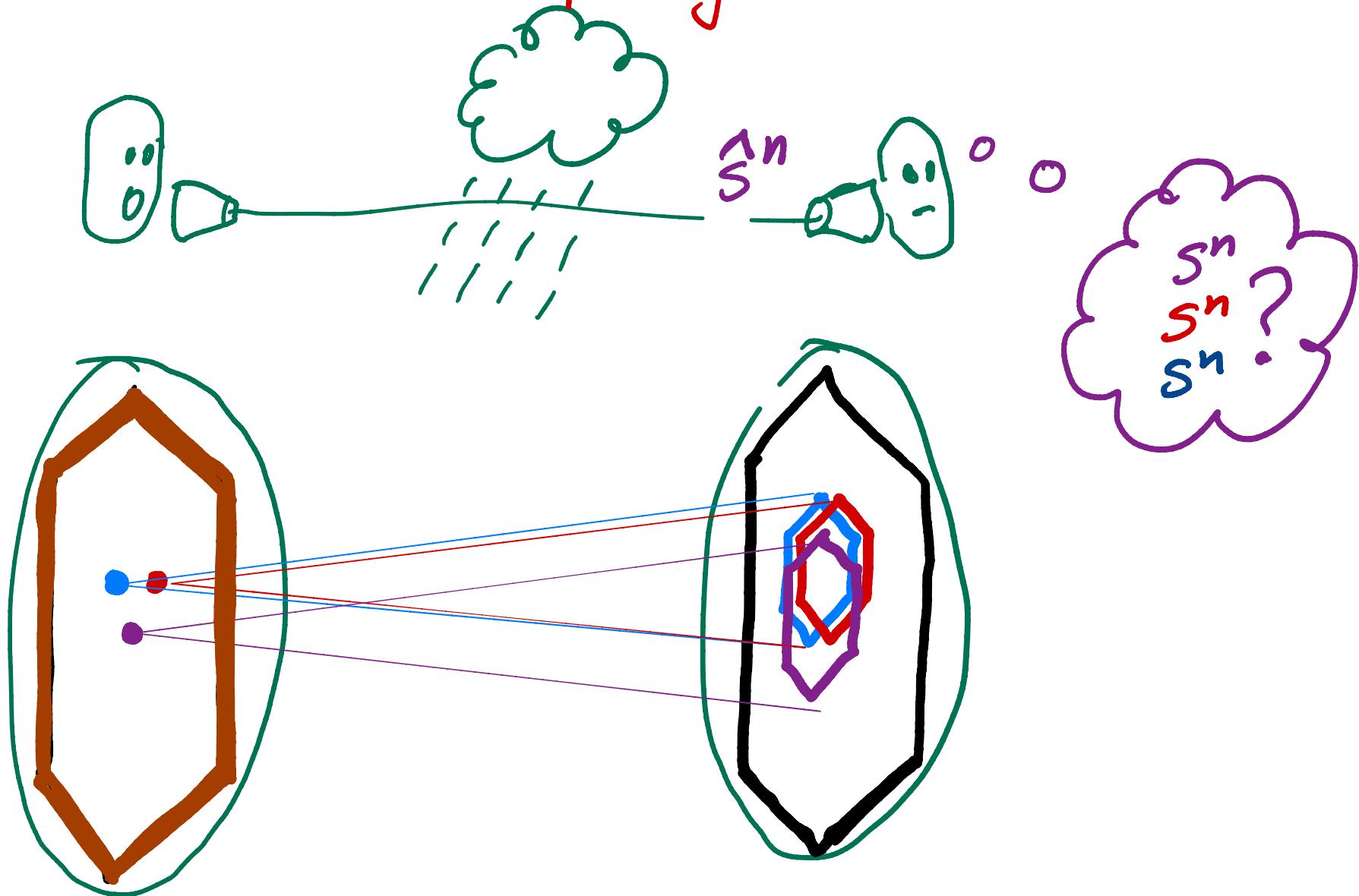
You can send more  
msgs!

What's the max # of msgs we can communicate?

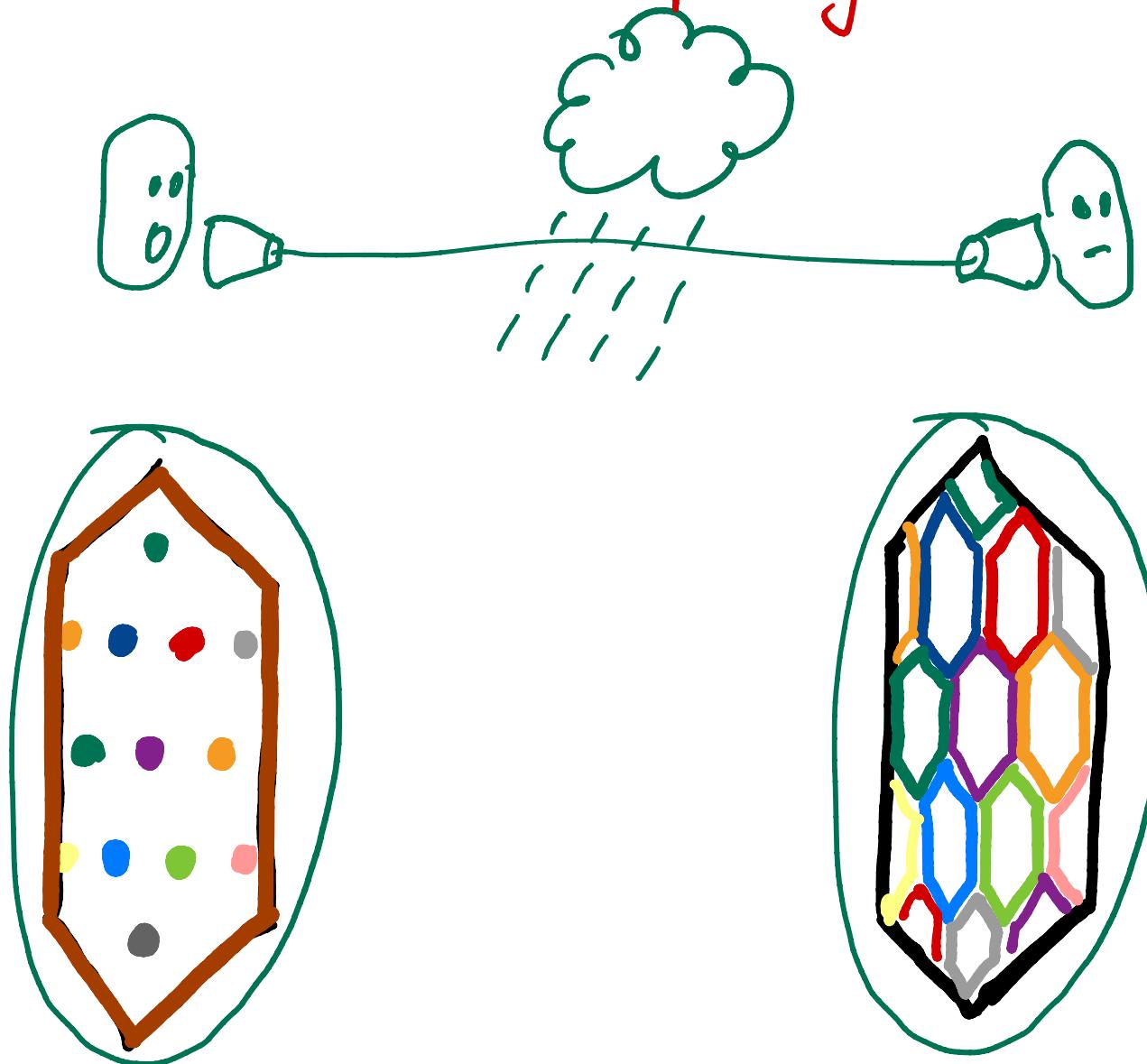


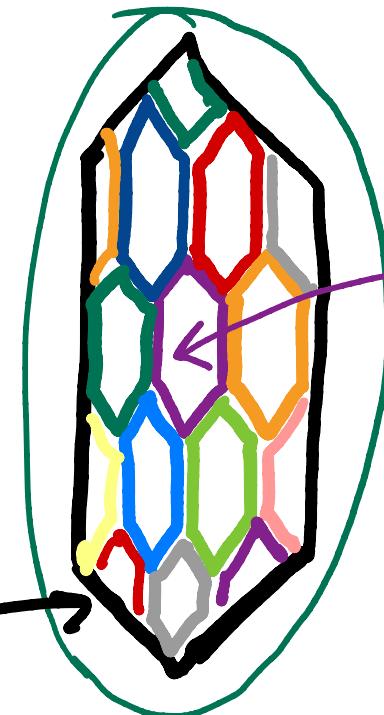
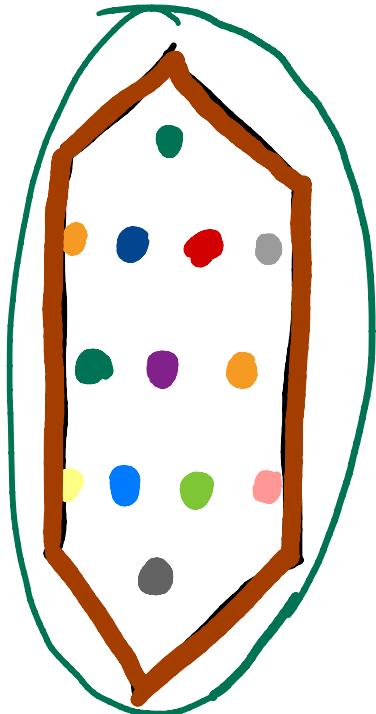
Too many messages!

What's the max # of msgs we can communicate?



What's the max # of msgs we can communicate?



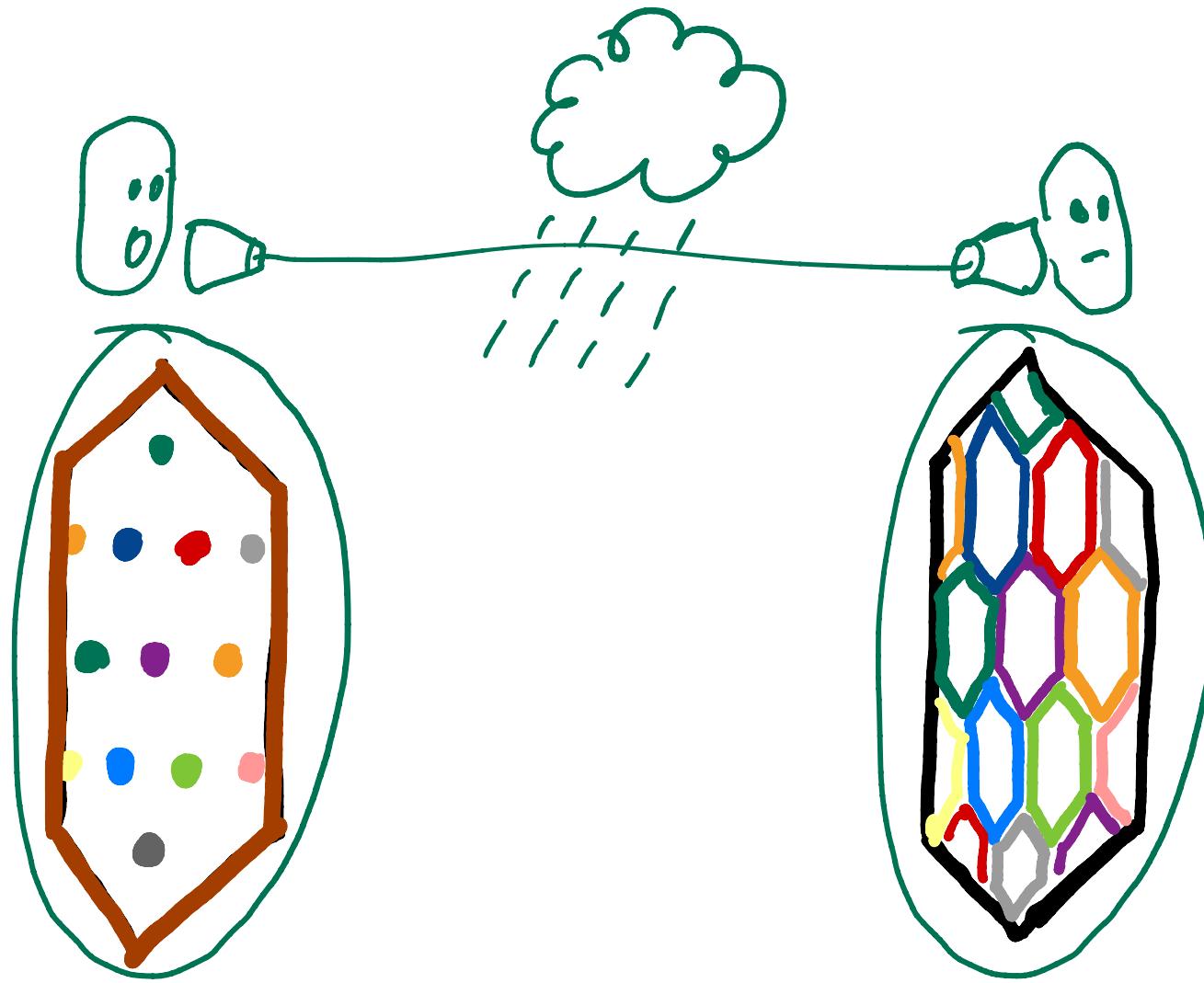


Typical set  
w.r.t.  $P_{\hat{S}} | S$   
 $\approx 2^{nH(\hat{S}|S)}$

Typical set  
w.r.t.  $P_{\hat{S}}$   $\approx 2^{nH(\hat{S})}$

$$\max \# \text{ of msgs} \leq \frac{2^{nH(\hat{S})}}{2^{nH(\hat{S}|S)}} = 2^{n(H(\hat{S}) - H(\hat{S}|S))}$$

$\stackrel{\text{def}}{=} 2^{nI(S; \hat{S})} = 2^n C$



If you choose  $2^{nr}$  codewords at random,  
then w.h.p., the typical sets  
don't overlap!

## Definition

\* Channel:  $(\mathcal{A}, \{p(\cdot | a)\}_{a \in \mathcal{A}})$

\* An  $(M, n)$ -code is a pair  $(\text{Enc}, \text{Dec})$

$$\text{Enc}: \{1, \dots, M\} \rightarrow \mathcal{A}^n$$

$$\text{Dec}: \mathcal{A}^n \rightarrow \{1, \dots, M\}$$

\* The rate is  $\frac{\log M}{n}$  bits/transmission

\* Prob. comm error =  $\frac{1}{M} \sum_{i \leq M} \Pr(\text{Dec}(\hat{S}^n) \neq i | S^n = \text{Enc}(i))$

## Typicality:

Thm: (Asymptotic equipartition) If  $S_1, S_2, \dots, S_n$  are i.i.d.  $\sim p(\cdot)$  then

$$-\frac{1}{n} \log p(\underbrace{S_1, \dots, S_n}_{\text{codeword}}) \xrightarrow[n \rightarrow \infty]{\parallel} H(p) \quad \text{in prob.}$$

As  $n \rightarrow \infty$ , if you choose a codeword  $s^n$  at random, by choosing each of its symbols at random using  $p(\cdot)$ , then  $p(s^n) \approx 2^{-nH(p)}$  w.h.p.

Proof:

$$-\frac{1}{n} \log p(s_1, \dots, s_n)$$

$$= -\frac{1}{n} \sum_{i \leq n} \log p(s_i) = \frac{1}{n} \sum_{i \leq n} -\log p(s_i)$$

$$\xrightarrow{n \rightarrow \infty} \mathbb{E}_S [-\log p(s)] \quad \begin{matrix} \text{in prob.} \\ \text{with } S \sim p(\cdot) \end{matrix}$$

$$= H(p) = H(S).$$

Def: The typical set  $A_{\varepsilon}^{(n)}$  w.r.t.  $p(\cdot)$  is the set of sequences such that

$(s_1, \dots, s_n) \in A_{\varepsilon}^{(n)}$  if

$$\left| -\frac{1}{n} \log p(s_1, \dots, s_n) - H(p) \right| \leq \varepsilon.$$

A codeword is typical if  
 $p(\text{codeword}) \approx 2^{-nH(p)}$

Thm:

$p(\text{typical codeword})$

1) If  $(s_1 \dots s_n) \in A_{\varepsilon}^{(n)}$  then

$$\approx 2^{-nH(p)}$$

$$2^{-n(H(p)+\varepsilon)} \leq p(s_1 \dots s_n) \leq 2^{-n(H(p)-\varepsilon)}$$

2)  $\Pr[A_{\varepsilon}^{(n)}] \geq 1-\varepsilon$ ,

for sufficiently large  $n$

If you choose a codeword at random then w.h.p. that codeword is typical.

3)  $(1-\varepsilon)2^{n(H(p)-\varepsilon)} \leq |A_{\varepsilon}^{(n)}| \leq 2^{n(H(p)+\varepsilon)}$

for large  $n$

$$|A_{\varepsilon}^{(n)}| \approx 2^{nH(p)}.$$

Proof:

1) If  $(s_1, \dots, s_n) \in A_{\varepsilon}^{(a)}$ ,

$$\left| -\frac{1}{n} \log p(s_1, \dots, s_n) - H(p) \right| \leq \varepsilon$$

$$H(p) - \varepsilon \leq -\frac{1}{n} \log p(s_1, \dots, s_n) \leq H(p) + \varepsilon$$

$$2^{-n(H(p) + \varepsilon)} \leq p(s_1, \dots, s_n) \leq 2^{-n(H(p) - \varepsilon)}$$

$$p(\text{typical codeword}) \approx 2^{-nH(p)}$$

Proof:

2)  $\Pr[A_{\varepsilon}^{(n)}] \geq 1 - \varepsilon$ , for sufficiently large  $n$ .

$$-\frac{1}{n} \log p(s_1, \dots, s_n) \xrightarrow[n \rightarrow \infty]{\text{in-prob}} H[p]$$

for  $\varepsilon > 0$ ,  
there is  $N$ ,  
for  $n \geq N$

$$\Pr\left(\left|-\frac{1}{n} \log p(s_1, \dots, s_n) - H[p]\right| < \varepsilon\right) \geq 1 - \varepsilon$$

$(s_1, \dots, s_n) \in A_{\varepsilon}^{(n)}$

$$\Pr((s_1, \dots, s_n) \in A_{\varepsilon}^{(n)}) \geq 1 - \varepsilon$$

$$\Pr(A_{\varepsilon}^{(n)}) \geq 1 - \varepsilon.$$

If you choose a codeword  $s^n$  at random, then  
w.h.p.  $s^n$  is typical w.r.t.  $p(\cdot)$ .

Proof

3.  $(1-\varepsilon)2^{n(H(p)-\varepsilon)} \leq |A_{\varepsilon}^{(n)}| \leq 2^{n(H(p)+\varepsilon)}$  for large  $n$

$$\begin{aligned} 1 &= \sum_{(s_1 \dots s_n)} p(s_1 \dots s_n) \geq \sum_{(s_1 \dots s_n) \in A_{\varepsilon}^{(n)}} p(s_1 \dots s_n) \\ &\geq \sum_{(s_1 \dots s_n) \in A_{\varepsilon}^{(n)}} 2^{-n(H(p)+\varepsilon)} \\ &= 2^{-n(H(p)+\varepsilon)} |A_{\varepsilon}^{(n)}| \end{aligned}$$

$$|A_{\varepsilon}^{(n)}| \leq 2^{-n(H(p)+\varepsilon)}$$

Proof

3.  $\underline{(1-\varepsilon)2^{n(H(p)-\varepsilon)} \leq |A_\varepsilon^{(n)}| \leq 2^{n(H(p)+\varepsilon)}} \text{ for large } n$

$$\Pr[A_\varepsilon^{(n)}] \geq 1-\varepsilon, \text{ for suff. large } n$$

$$1-\varepsilon \leq \Pr[A_\varepsilon^{(n)}] = \sum_{(s_1 \dots s_n) \in A_\varepsilon^{(n)}} p(s_1 \dots s_n) \leq \sum_{(s_1 \dots s_n) \in A_\varepsilon^{(n)}} 2^{-n(H(p)-\varepsilon)}$$

$$\Rightarrow (1-\varepsilon)2^{n(H(p)-\varepsilon)} \leq |A_\varepsilon^{(n)}| = 2^{-n(H(p)-\varepsilon)} |A_\varepsilon^{(n)}|$$

$$|A_\varepsilon^{(n)}| \approx 2^{nH(p)}$$

Def: The set  $A_{\varepsilon}^{(n)}$  of jointly typical pairs of seqs w.r.t.  
 $P_{S,\hat{S}}(\cdot, \cdot)$  is the set of pairs of sequences  
such that

$$(s^n, \hat{s}^n) \in A_{\varepsilon}^{(n)} \text{ iff}$$

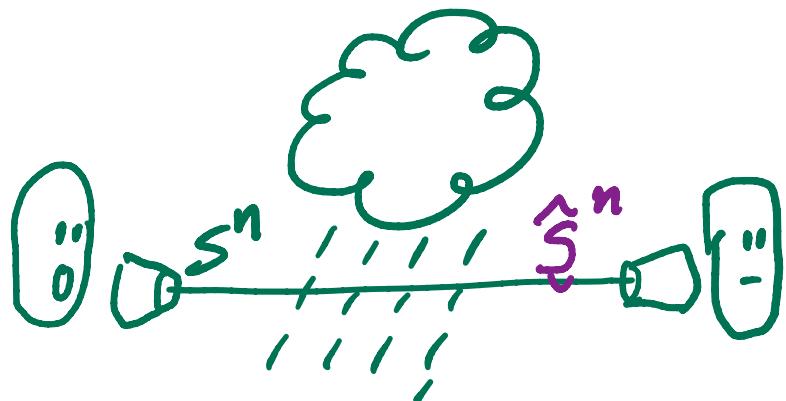
$$1. \left| -\frac{1}{n} \log P_{S^n}(s^n) - H(S) \right| < \varepsilon \quad 2. \left| -\frac{1}{n} \log P_{\hat{S}^n}(\hat{s}^n) - H(\hat{S}) \right| < \varepsilon$$

$$3. \left| -\frac{1}{n} \log P_{S^n, \hat{S}^n}(s^n, \hat{s}^n) - H(S, \hat{S}) \right| < \varepsilon$$

In our context :  $S^n$  input codeword,  $\hat{S}^n$  output codeword

Thm :

1.  $\Pr((s^n, \hat{s}^n) \in A_{\varepsilon}^{(n)}) \rightarrow 1$   
as  $n \rightarrow \infty$

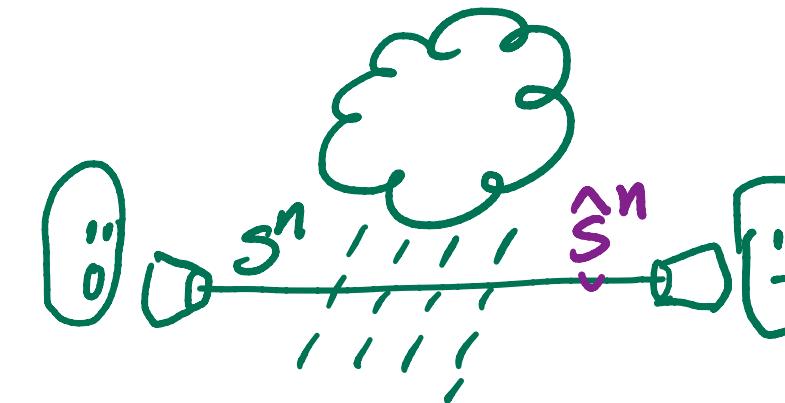


If  $s^n$  drawn from  $P_{S^n}^*$   
then  $(s^n, \hat{s}^n) \in A_{\varepsilon}^{(n)}$   
w.h.p.

If  
1.  $s^n \xleftarrow{\$} P_{S^n}^*(\cdot)$   
2.  $\hat{s}^n \xleftarrow{\$} P_{\hat{S}^n | S^n = s^n}(\cdot)$   
then  $(s^n, \hat{s}^n) \in A_{\varepsilon}^{(n)}$  w.h.p.

Thm :

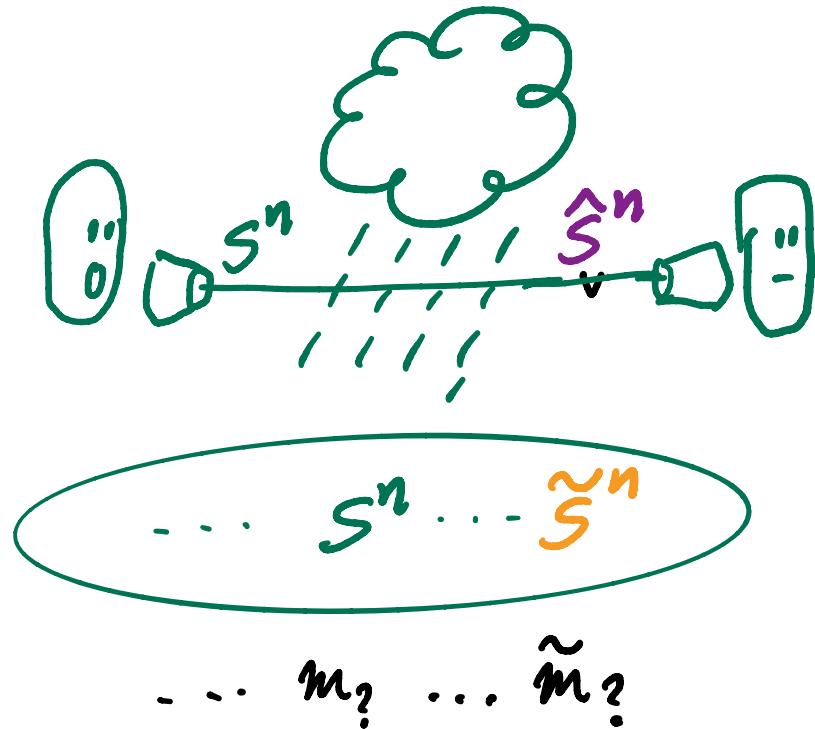
2. If  $\tilde{S}^n \sim p_{S^n}(\cdot)$  and  $\hat{S}^n \sim \hat{p}_{\hat{S}^n}(\cdot)$ . then  
 $\Pr[\tilde{S}^n, \hat{S}^n \in A_\epsilon^{(n)}] \leq 2^{-n}(I(S; \hat{S}) - 3\epsilon)$



$\dots S^n \dots \tilde{S}^n$

$\dots m_? \dots \tilde{m}_?$

$(\tilde{S}^n, \hat{S}^n) \notin A_\epsilon^{(n)}$  w.h.p.



$(s^n, \hat{s}^n) \in A_{\varepsilon}^{(n)}$  and  $(\tilde{s}^n, \hat{s}^n) \notin A_{\varepsilon}^{(n)}$  w.h.p.

Shannon's coding thm

If  $r <$  channel's capacity then  
there is a code s.t.

\* rate =  $r$  (i.e., #msgs =  $2^{nr}$ )

\* prob. error  $\xrightarrow{n \rightarrow \infty} 0$

(The converse is also true)

If #msgs =  $2^{nr}$  with  $r < C$ , then there  
is a  $(L2^{nr}, n)$ -code with pr.error  $\xrightarrow{n \rightarrow \infty} 0$ .

Remember:

A  $(\lfloor 2^{nr} \rfloor, n)$ -code is a pair  $(\text{Enc}, \text{Dec})$

\* Enc:  $\{1 \dots \lfloor 2^{nr} \rfloor\} \rightarrow \mathcal{U}^n$

\* Dec:  $\mathcal{U}^n \rightarrow \{1 \dots \lfloor 2^{nr} \rfloor\}$

Prob. comm. error =  $\frac{1}{\lfloor 2^{nr} \rfloor} \sum_{m \in \{1 \dots \lfloor 2^{nr} \rfloor\}} \Pr(\text{Dec}(\hat{S}^n) \neq m | S^n = \text{Enc}(m))$

## Shannon's random code



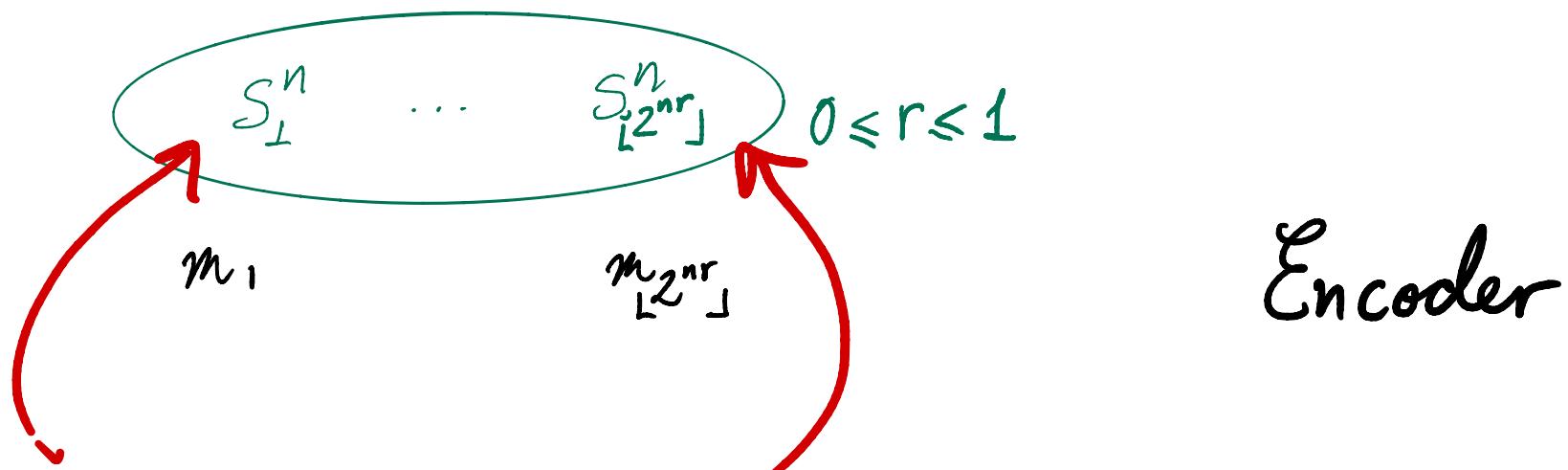
Let  $r <$  channel cap.

$m_1$

$m_{\lfloor 2^{nr} \rfloor}$

Pick  $\lfloor 2^{nr} \rfloor$  msgs.

# Shannon's random code

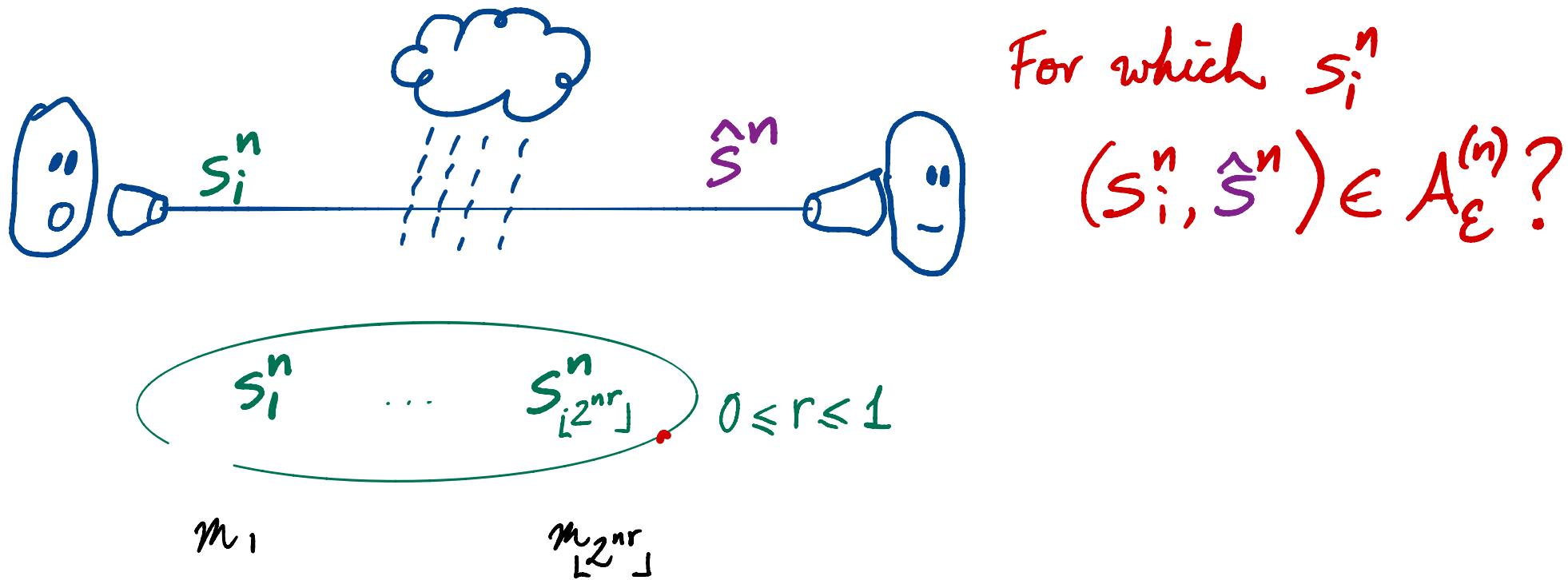


Choose the codewords

at random from  $P_S^* = \operatorname{argmax}_P I(S; \hat{S})$

The codewords are typical w.h.p.

# Shannon's random code



Decoder

# Shannon's random code



For which  $s_i^n$   
 $(s_i^n, \hat{s}^n) \in A_{\epsilon}^{(n)}$ ?



$m_1$

$m_{2^{nr}}$

We now prove that

as  $n \rightarrow \infty$

\* rate =  $r$

\* prob. error  $\rightarrow 0$

$$\text{rate} = \frac{\log |2^{nr}|}{n} \xrightarrow{n \rightarrow \infty} r$$

Proof:

$$\Pr(\mathcal{E}) = \Pr(\mathcal{E} | \#1) \quad S^n = \text{Enc}(1)$$

$$= \Pr \left[ \begin{array}{l} \text{Enc}(1), \hat{S}^n \notin A_{\mathcal{E}}^{(n)} \text{ or} \\ \text{Enc}(2), \hat{S}^n \in A_{\mathcal{E}}^{(n)} \text{ or} \\ \vdots \\ \text{Enc}(\lfloor 2^n \rfloor), \hat{S}^n \in A_{\mathcal{E}}^{(n)} \end{array} \middle| \#1 \right]$$

$$\leq \Pr[\text{Enc}(1), \hat{S}^n \notin A_{\mathcal{E}}^{(n)} | \#1] + \sum_{2 \leq w \leq \lfloor 2^n \rfloor} \Pr[\text{Enc}(w), \hat{S}^n \in A_{\mathcal{E}}^{(n)} | \#1]$$

$$\leq \Pr\left[\text{Enc}(1), \hat{S}^n \notin A_{\varepsilon}^{(n)} | \#1\right] + \sum_{2 \leq w \leq \lfloor 2^{nr} \rfloor} \Pr\left[\text{Enc}(w), \hat{S}^n \in A_{\varepsilon}^{(n)} | \#1\right]$$

$\eta_n \xrightarrow[n \rightarrow \infty]{} 0$

$\leq 2^{-n(I(S; \hat{S}) - 3\varepsilon)}$

$$\leq \eta_n + \sum_{2 \leq w \leq \lfloor 2^{nr} \rfloor} 2^{-n(I(S; \hat{S}) - 3\varepsilon)}$$

$$= \eta_n + (\lfloor 2^{nr} \rfloor - 1) 2^{-n(I(S; \hat{S}) - 3\varepsilon)}$$

$$\leq \eta_n + 2^{nr} 2^{-n(I(S; \hat{S}) - 3\varepsilon)}$$

$$= \eta_n + 2^{-n(I(S; \hat{S}) - r - 3\varepsilon)} \quad \text{if } r < I(S; \hat{S}) \\ = C$$

$\xrightarrow{n \rightarrow \infty} 0$



Shannon's coding thm

If # of msgs =  $2^{nr}$  (with  $r < C$ )

$\Rightarrow$  there is a code

$$\text{prob. error} \xrightarrow[n \rightarrow \infty]{} 0$$

(The converse is true)

This means  $\Rightarrow$   $\max \# \text{msgs} = 2^{nC}$

$\Rightarrow$  capacity as a channel measure