

Validation via information theory II

Statistical Learning Theory 2020

Organization

What is PA?

Roadmap

Formalization of PA

→ Algorithms as channels

→ Derivation of PA

Applications

Posterior agreement

How to validate algorithms?

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

Expected log posterior agreement:

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

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:

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

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

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

Organization

What is PA?

Roadmap

Formalization of PA

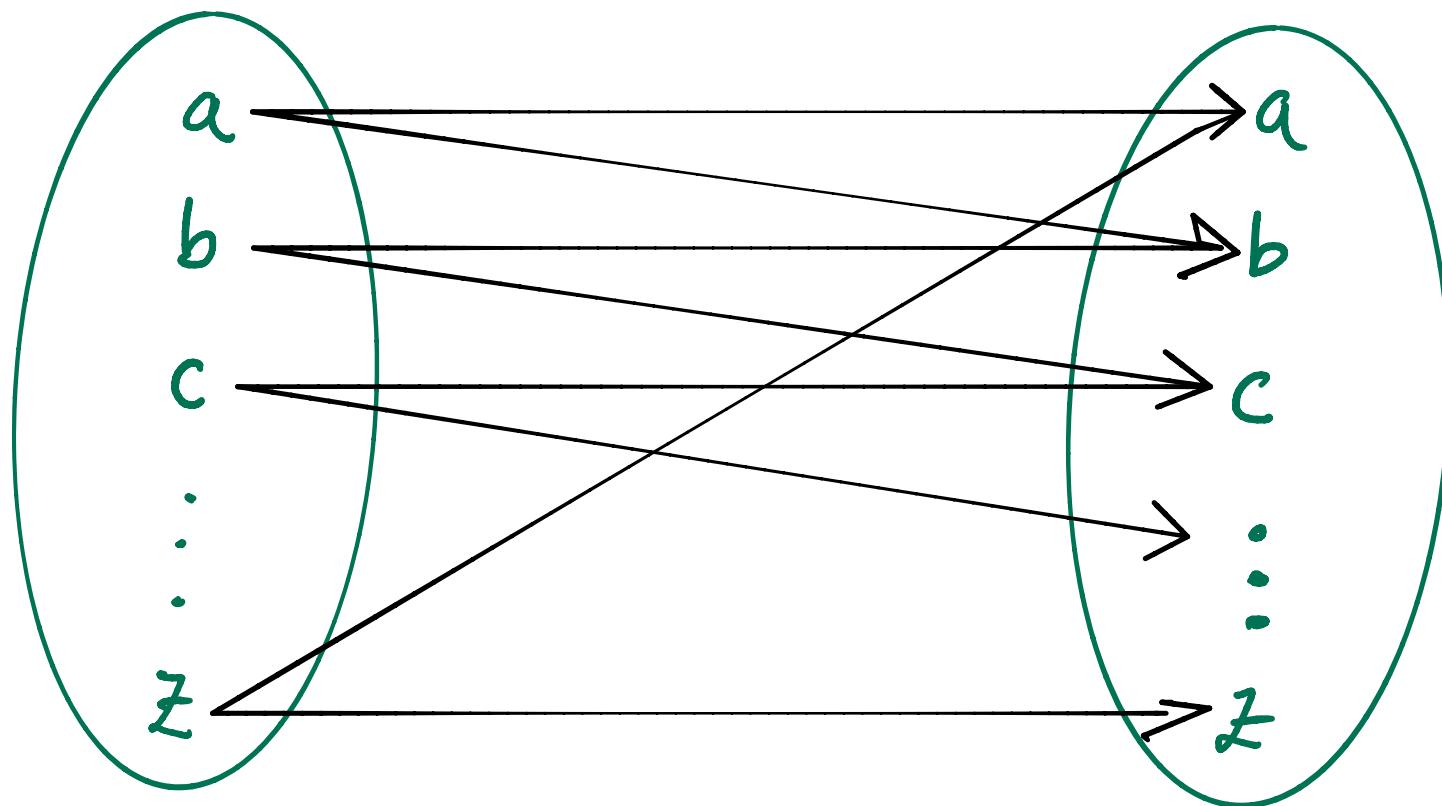
→ Algorithms as channels

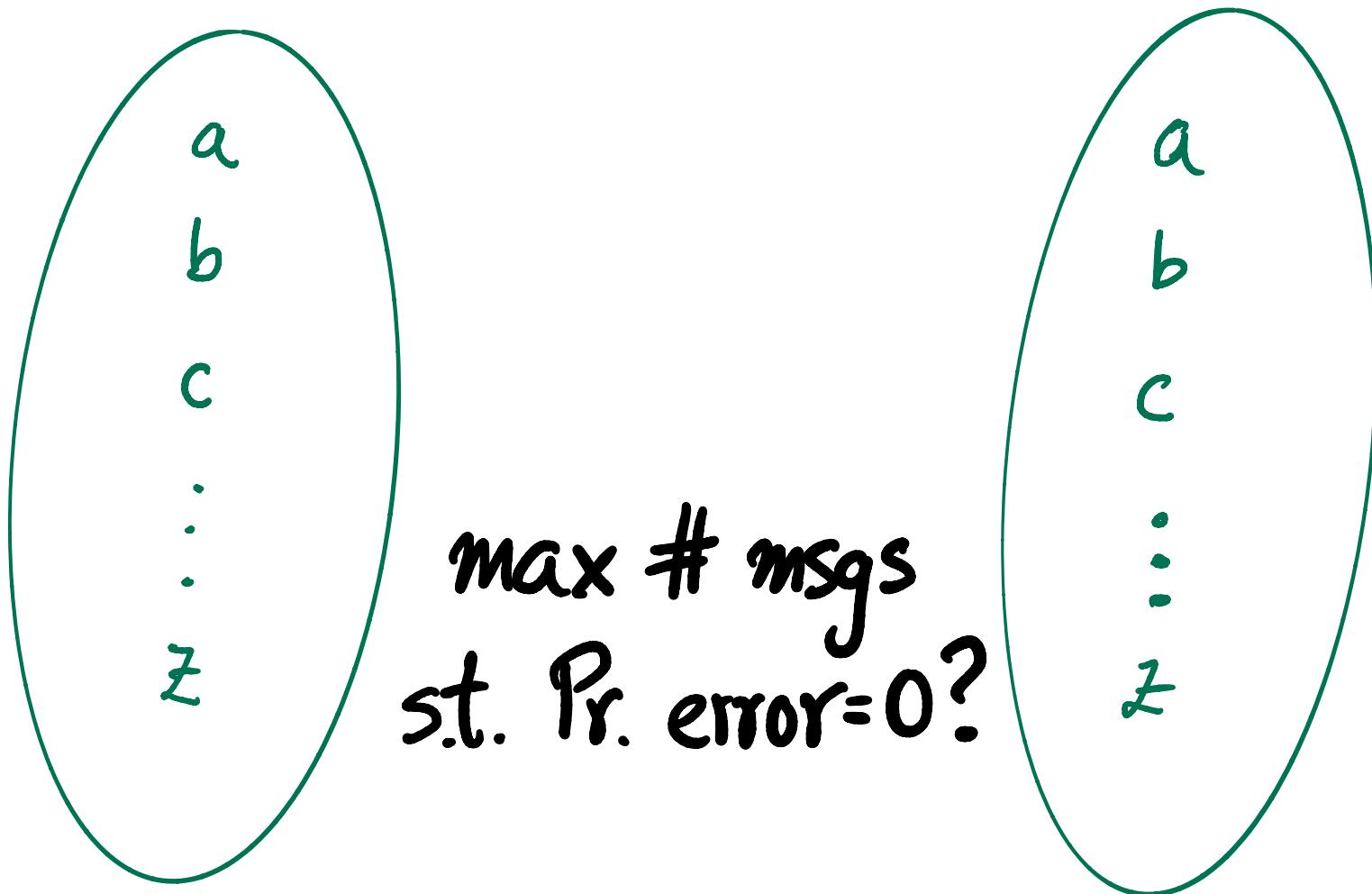
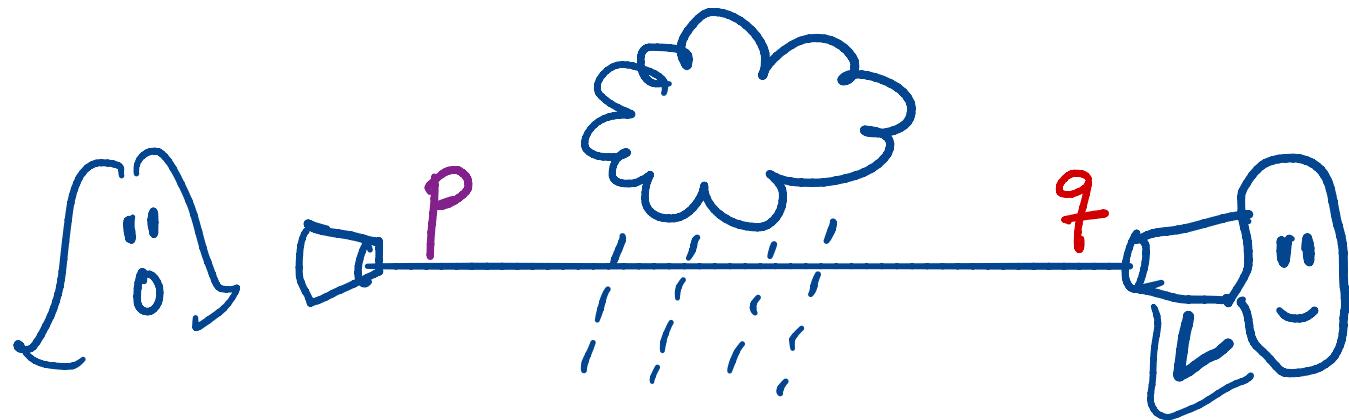
→ Derivation of PA

Applications

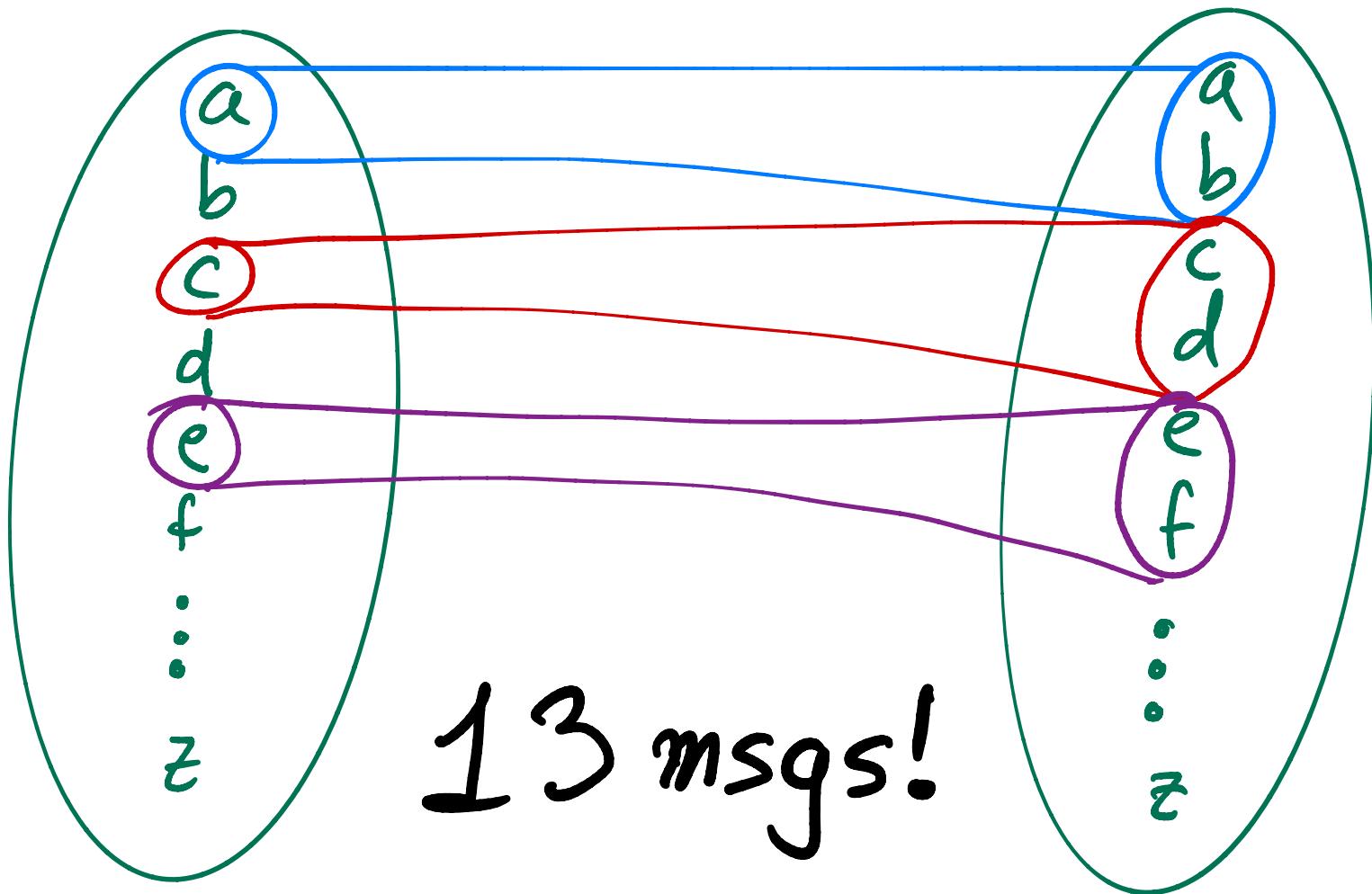
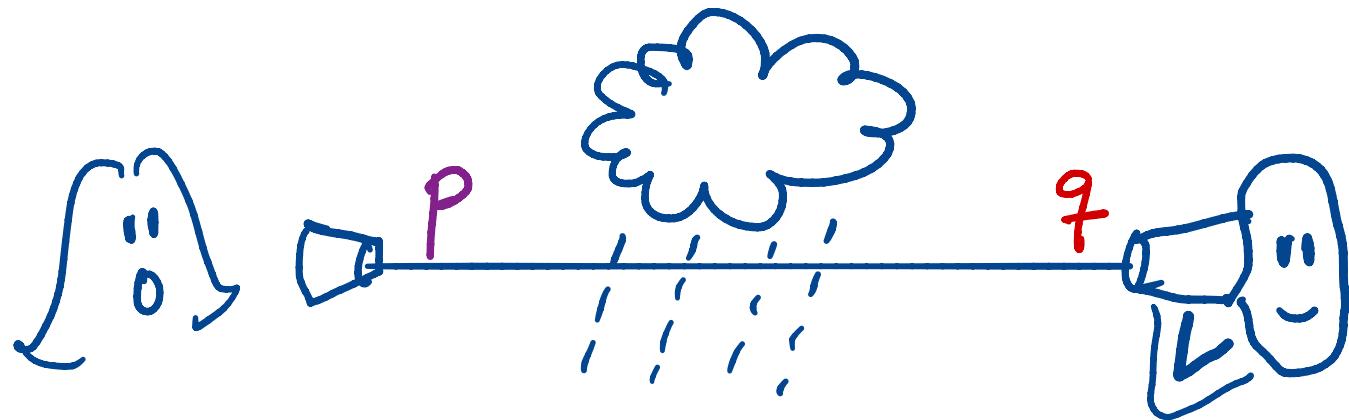
PA originates from
modeling and evaluating
algorithms as channels

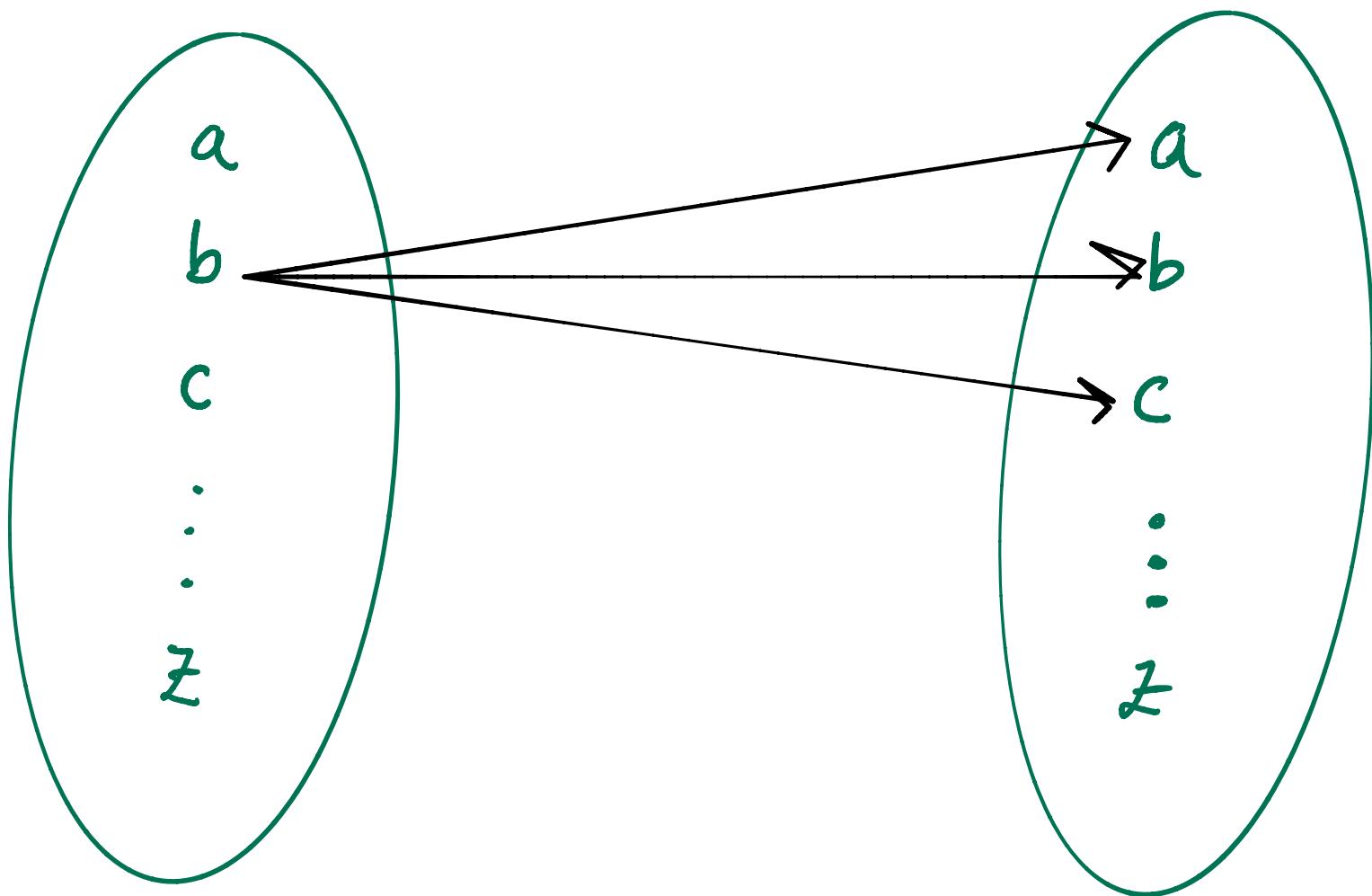
Communication on a noisy channel

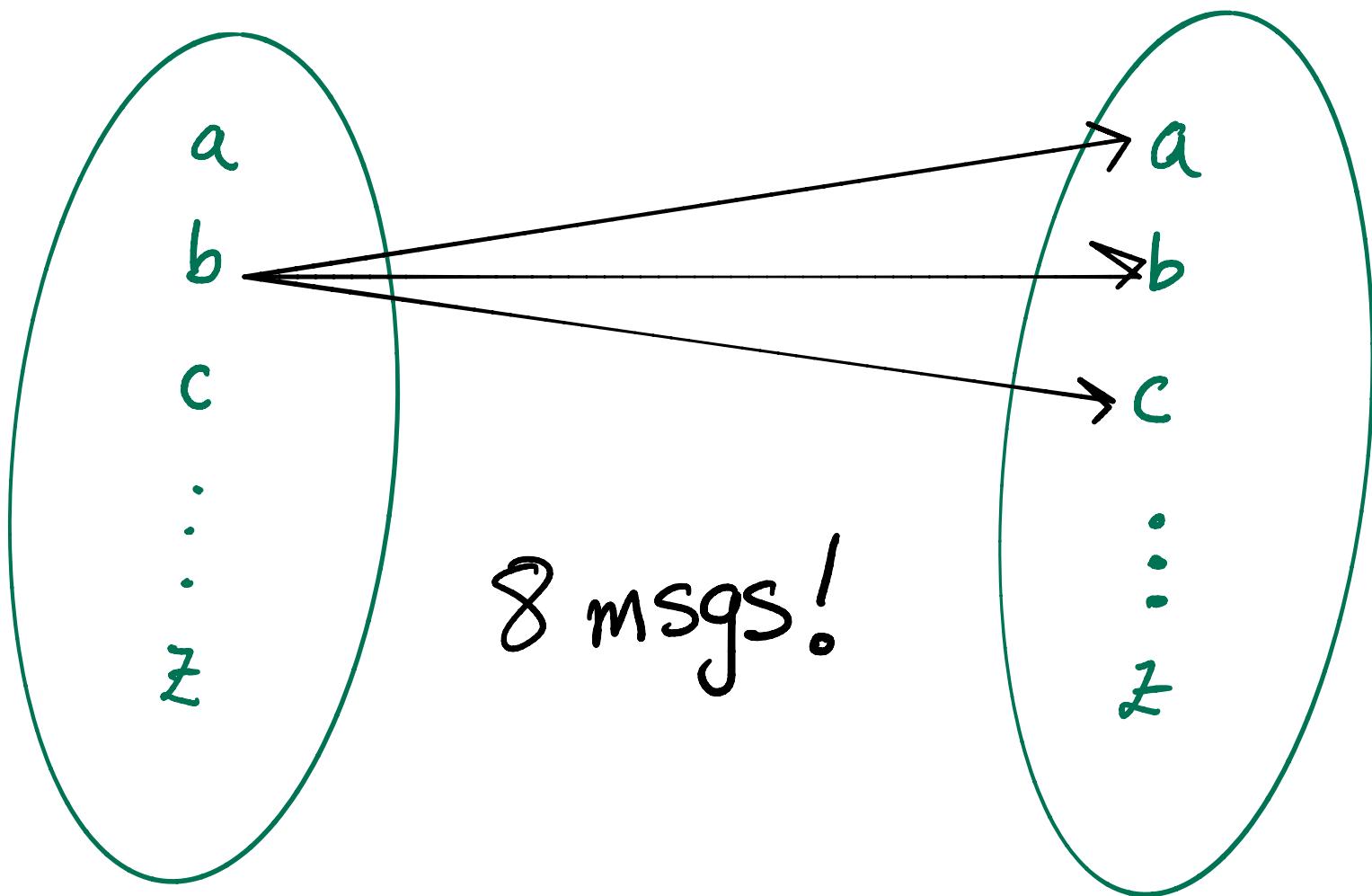


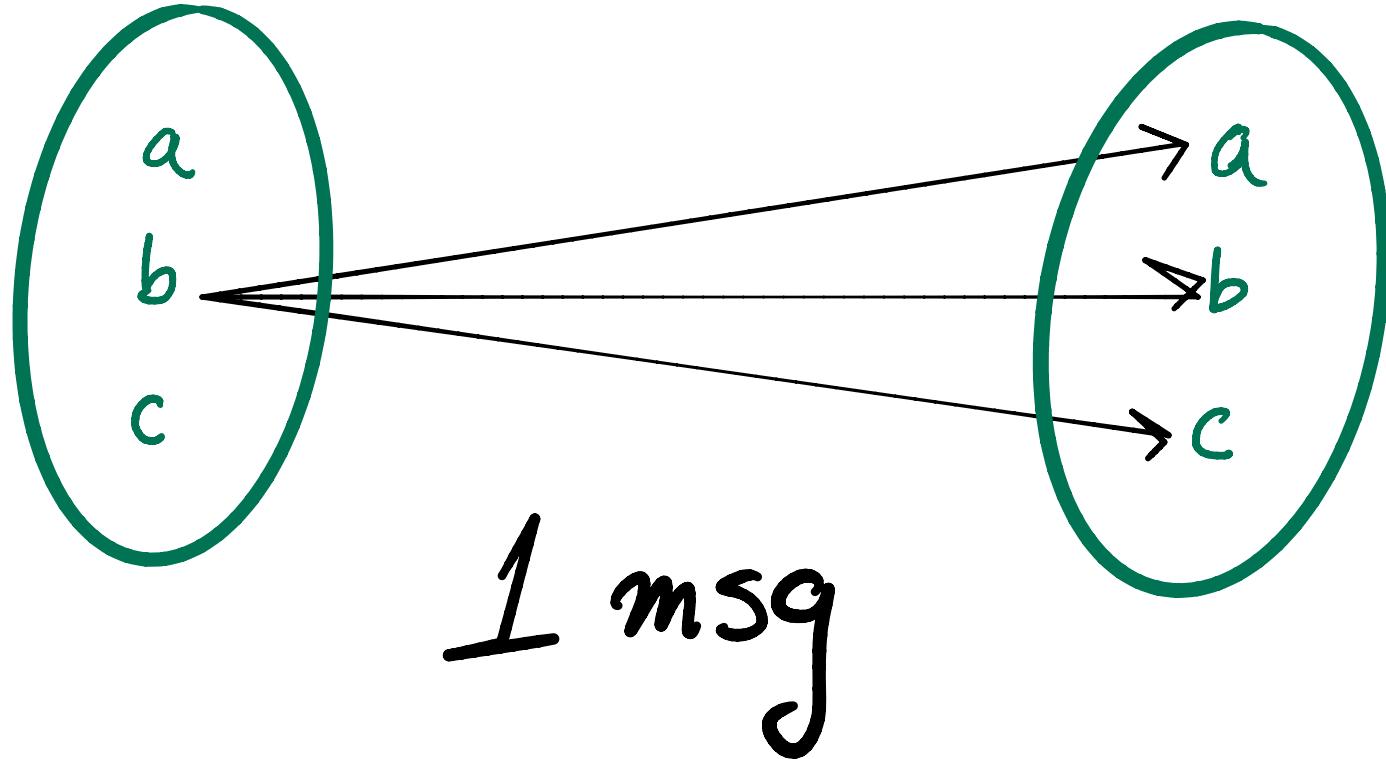


max # msgs
s.t. $\Pr.$ error=0?









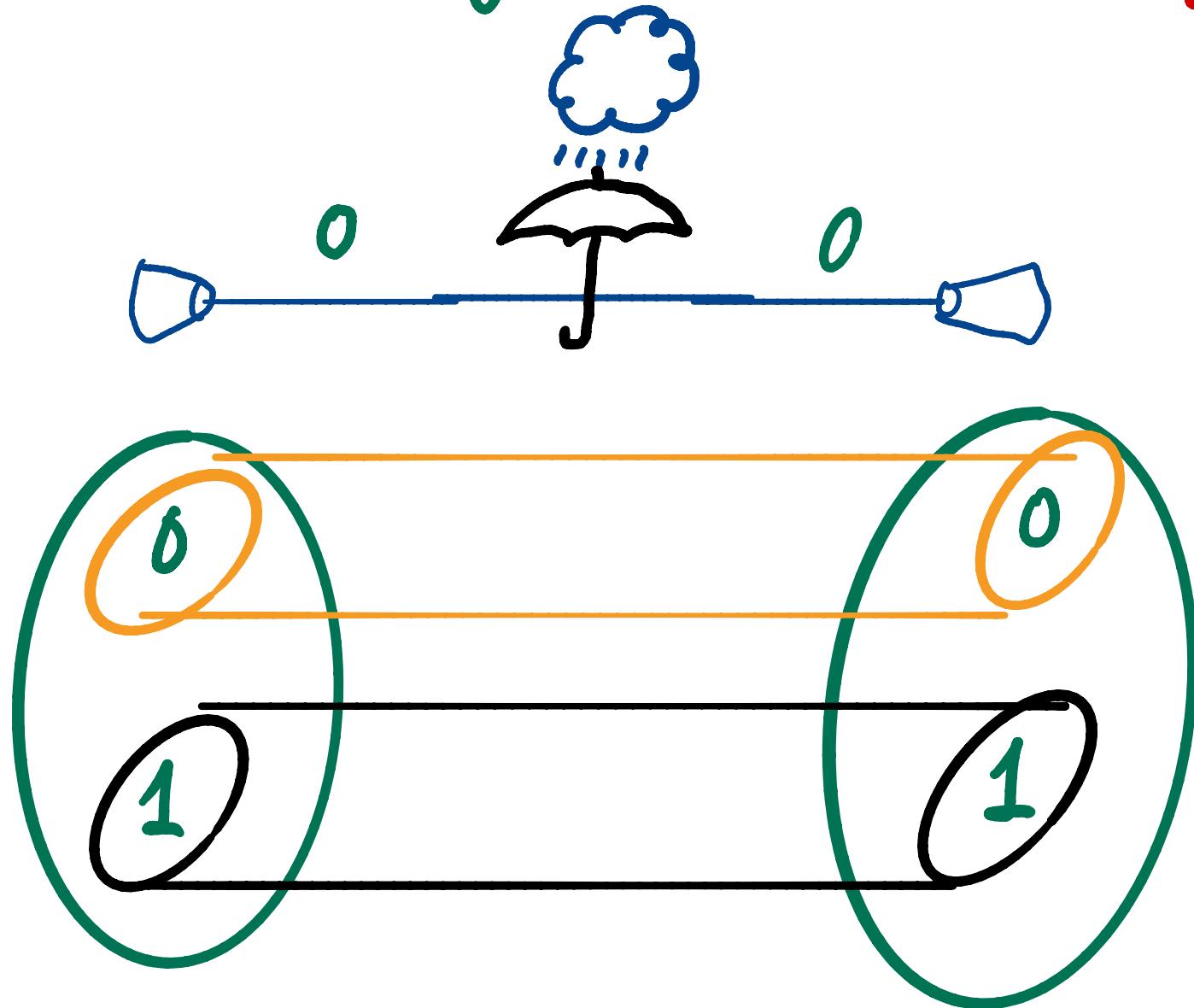
Good channels

\Leftrightarrow

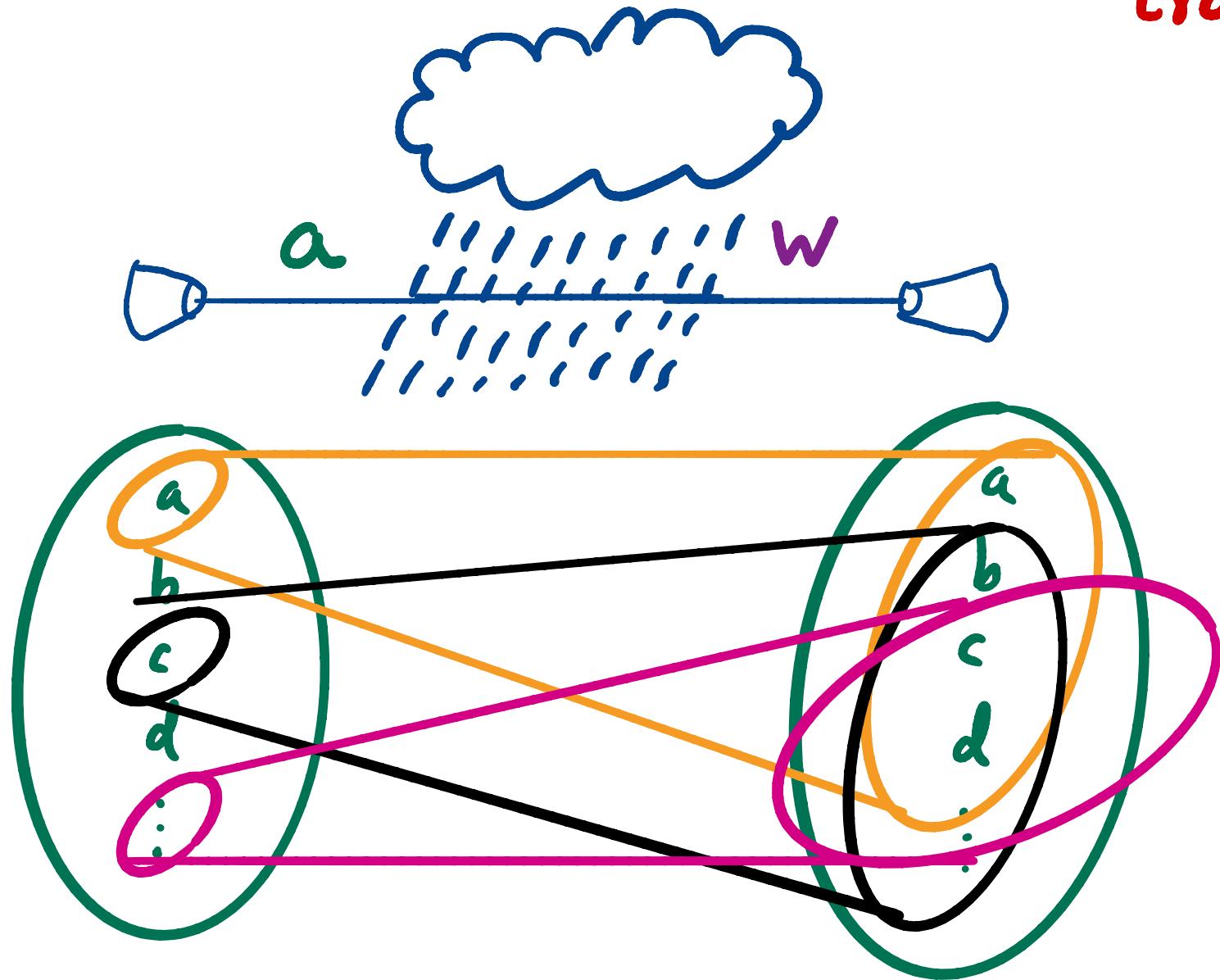
high info transfer

= codewords + channel
have robust
content against noise

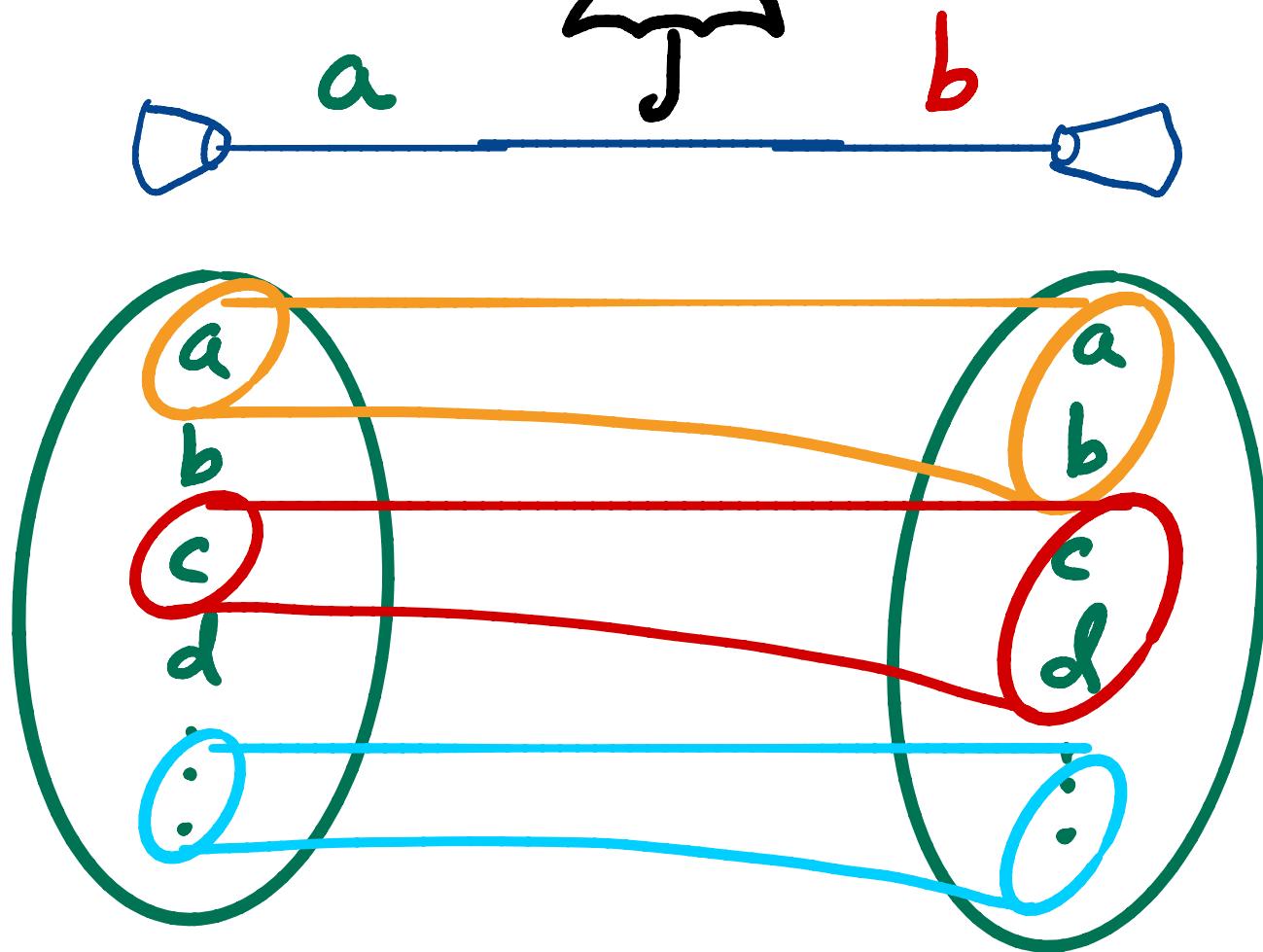
*poor
Content* + *robust
against noise* = *low
information
transfer*



rich
content + sensitive
to noise = low
information
transfer



rich
content + robust
against
noise = high
information
transfer



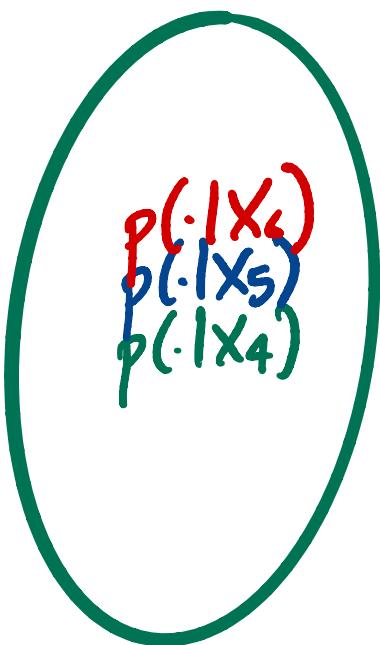
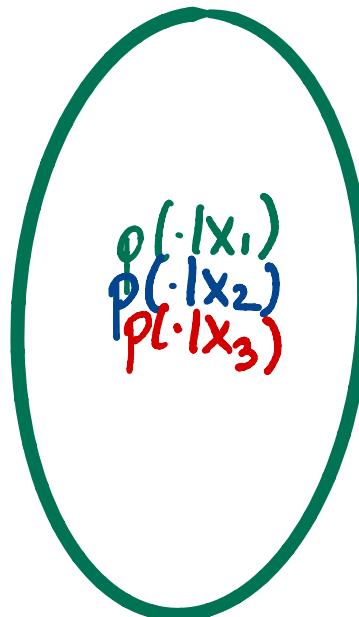
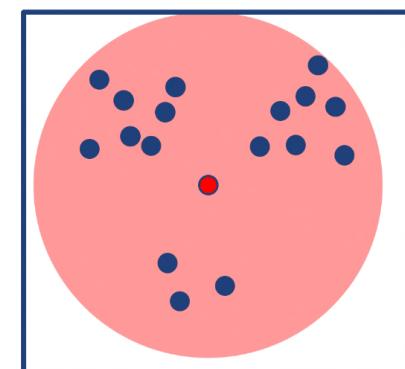
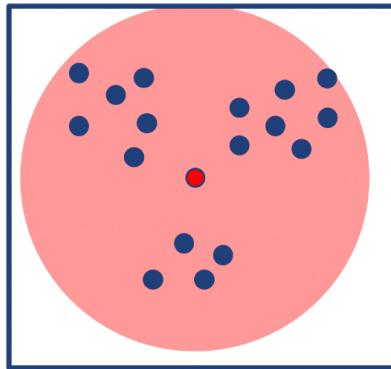
Good algorithms



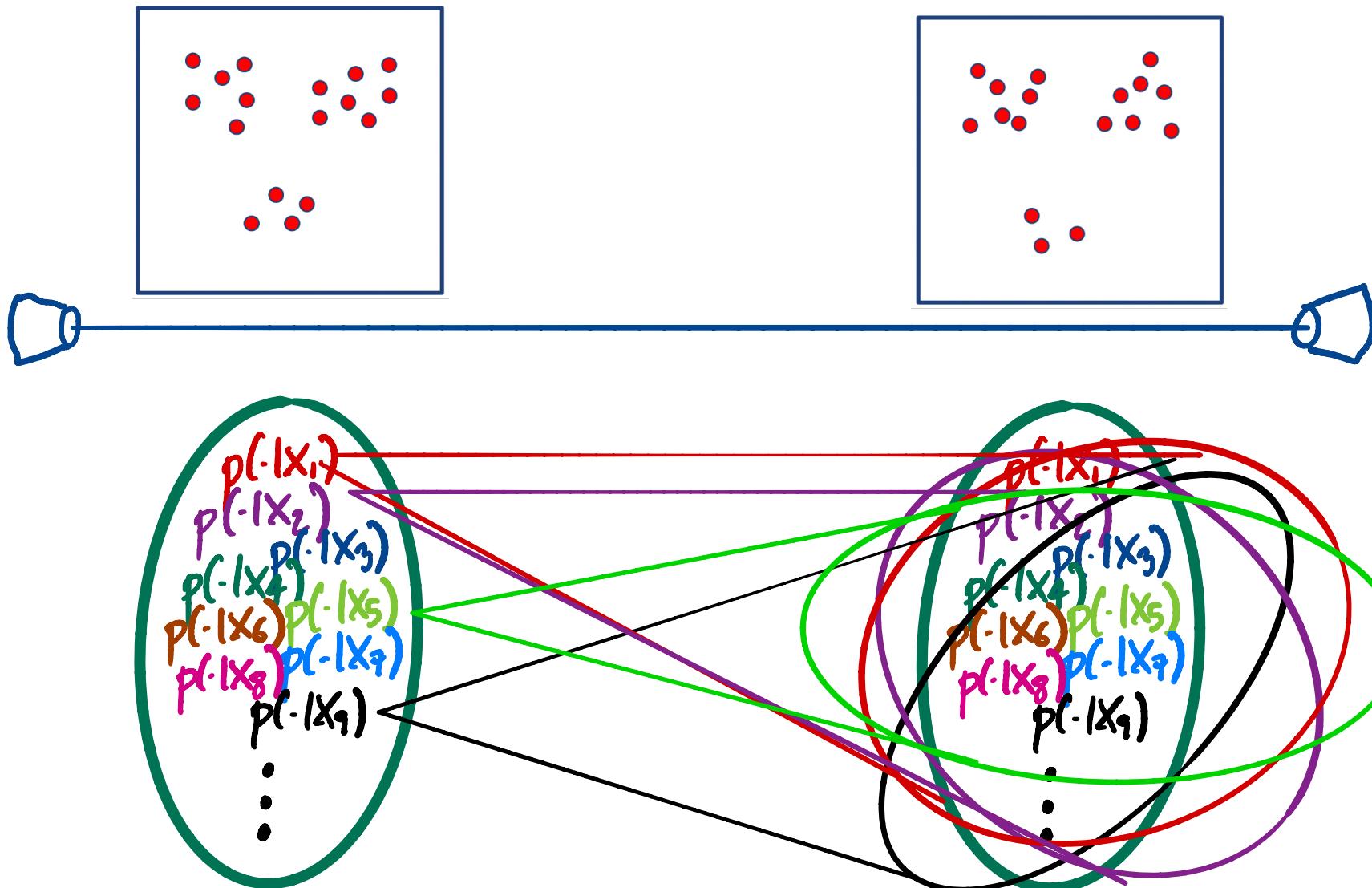
high info extraction

= outputs + algorithm
 have robust
 content against noise

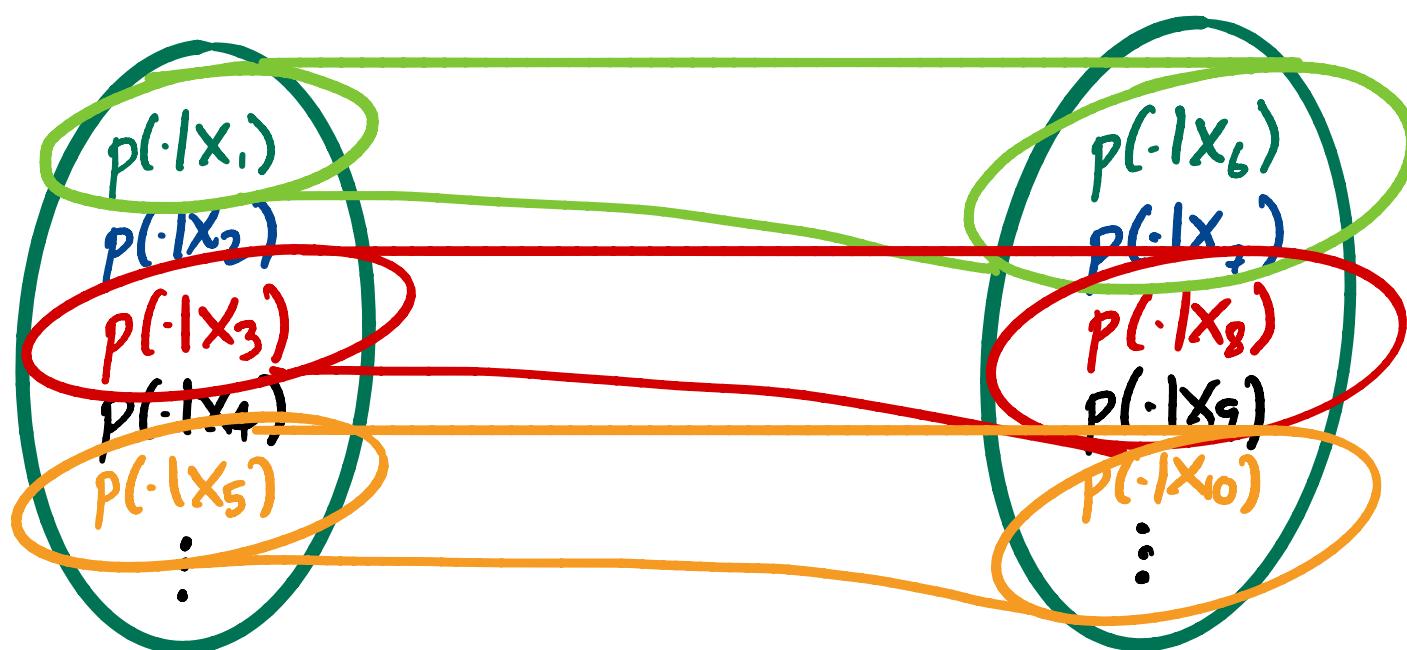
*poor
content* + *robust* = *low
info
extraction*



rich content + no robustness = low info extraction



rich content + robust = high info extraction



Information extraction

* Outputs have content

* Algorithm is robust

Adjust your algos
to increase exp. log. PA

Information transmission

* Codewords have content

* Channel is robust

Adjust your channels
to increase 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?

Roadmap

Formalization of PA

→ Algorithms as channels

→ Derivation of PA

Applications

Algorithms as
communication channels

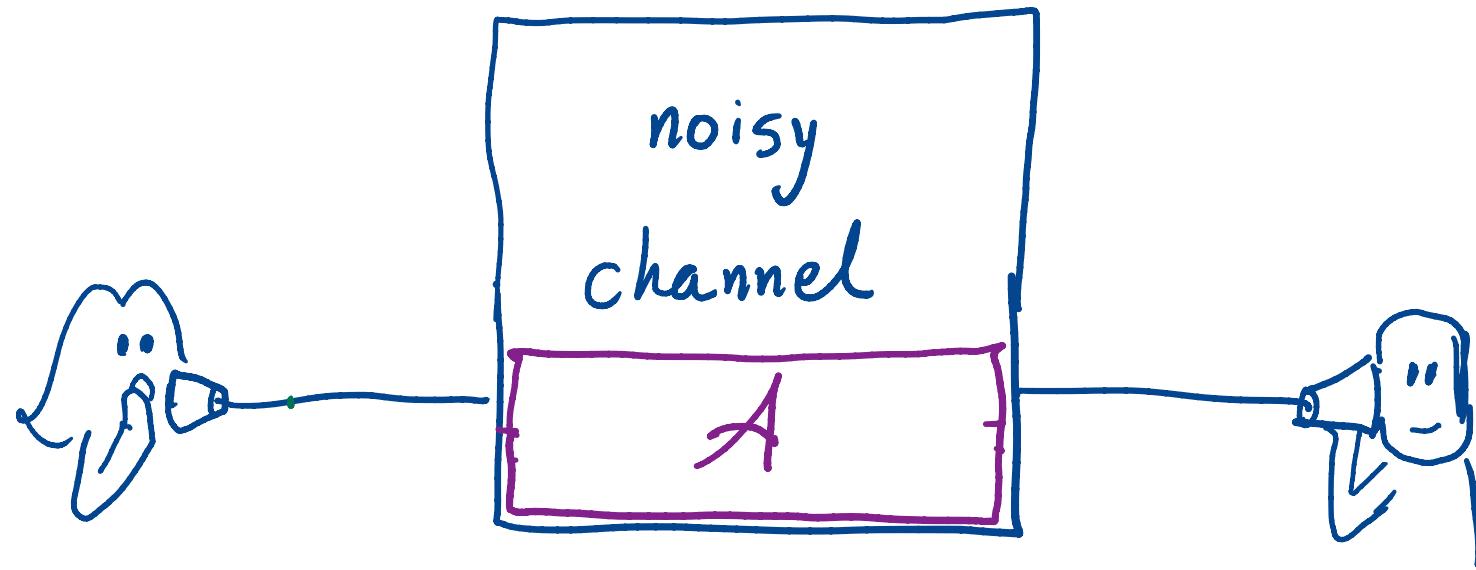
Assumptions

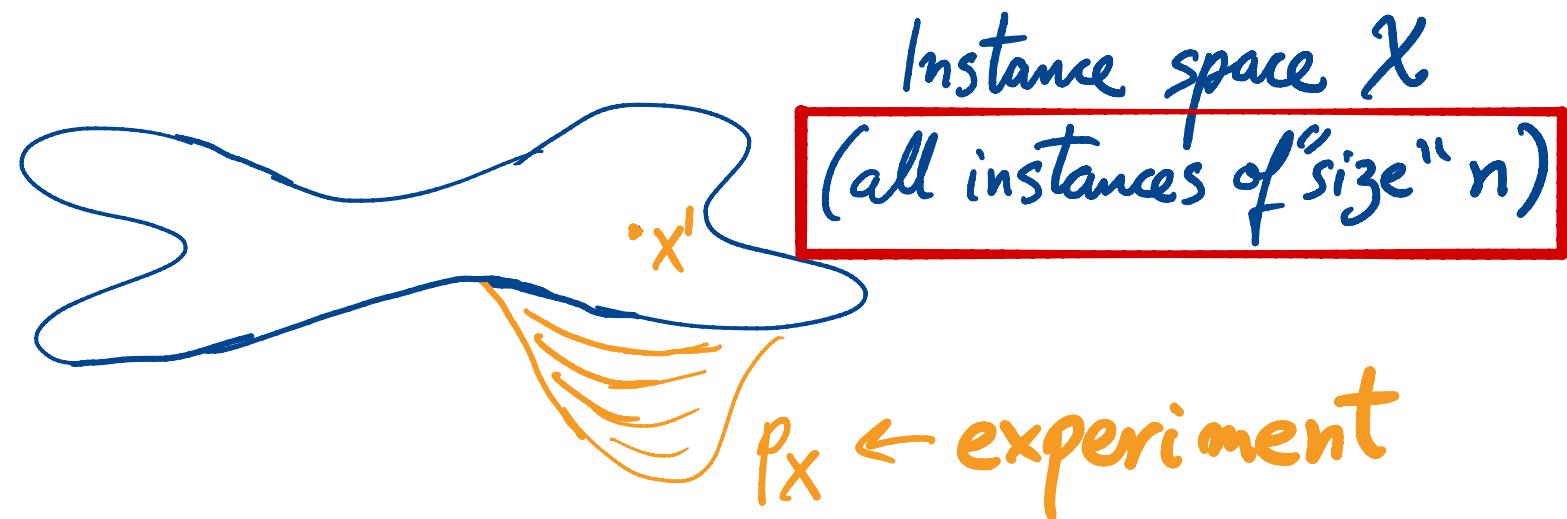
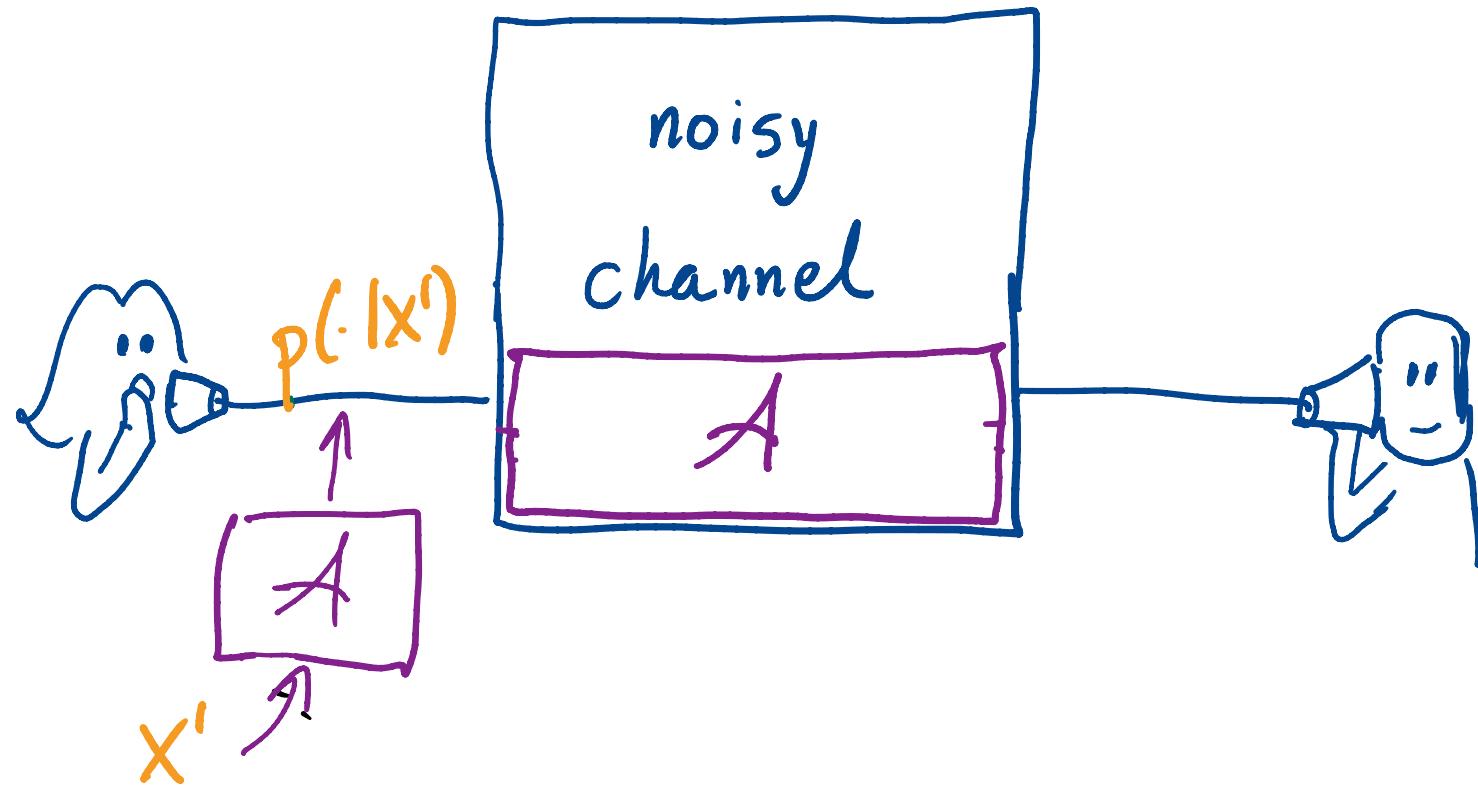
* Exponential sol space: $\log |G| = O(n)$.

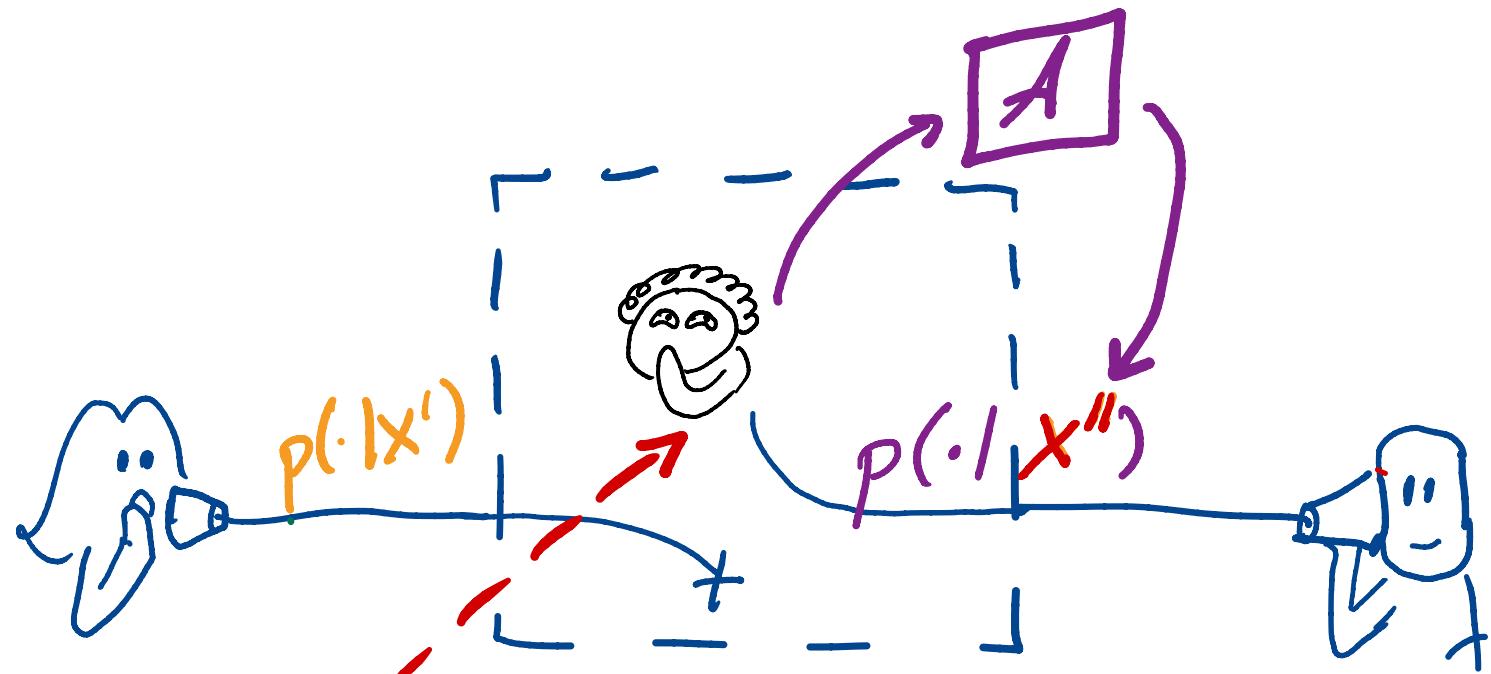
* Prob outputs: $X' \xrightarrow{\text{A}} p(\cdot | X')$

* Asymptotic equip. property:

$$\underbrace{\text{exp. log. PA}}_{\approx \frac{1}{\log |C|} \sum_i \log y_i} \xrightarrow[n \rightarrow \infty]{\text{in prob}} \underbrace{\text{exp. log. PA}}_{\approx \mathbb{E}_y[\log y]}$$



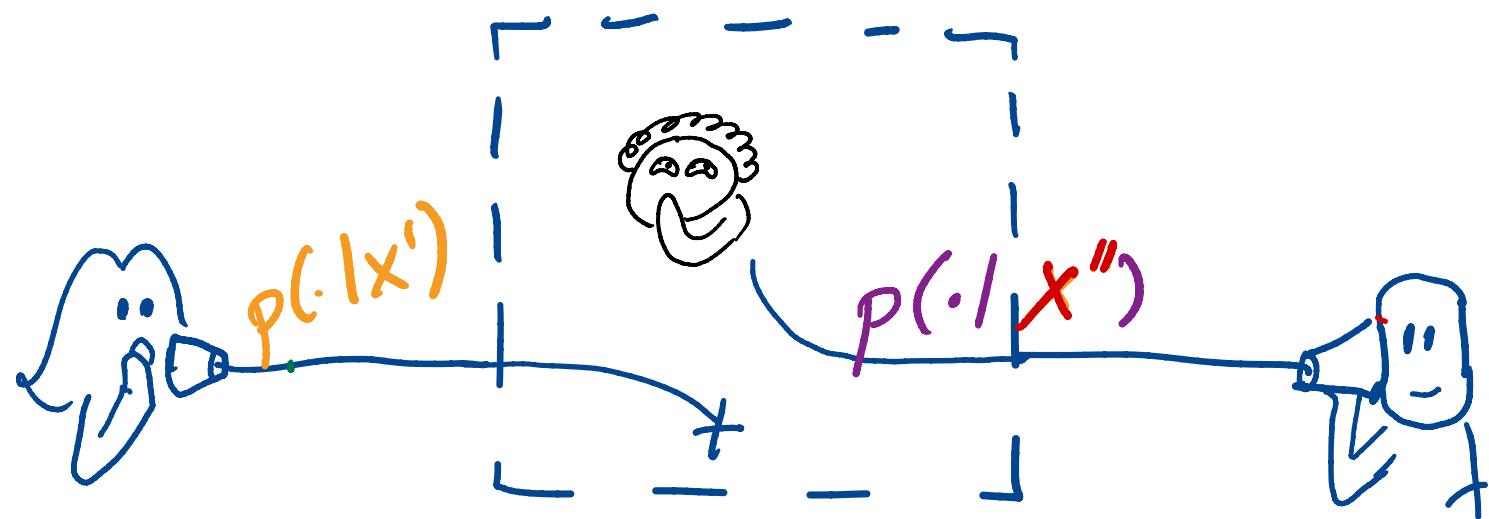




x''

Instance space X

p_X

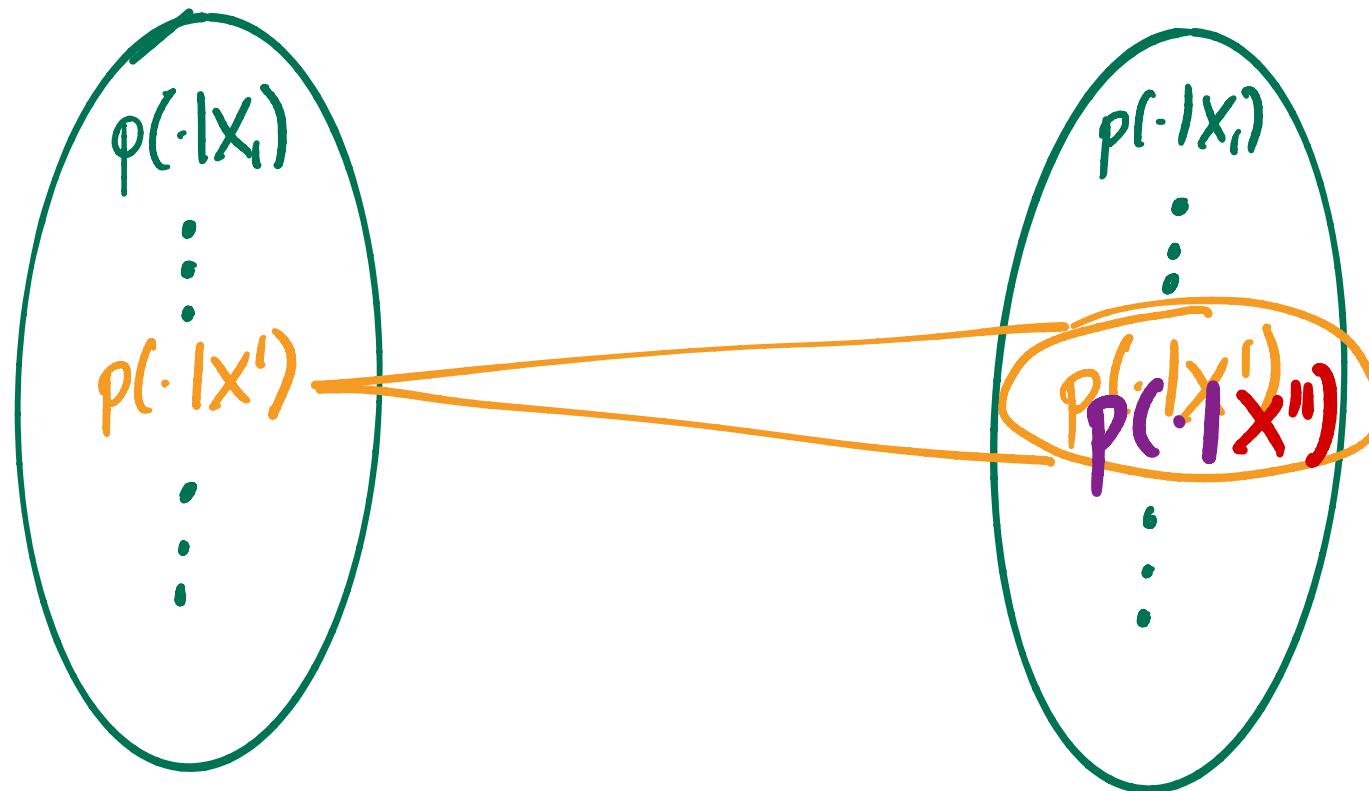
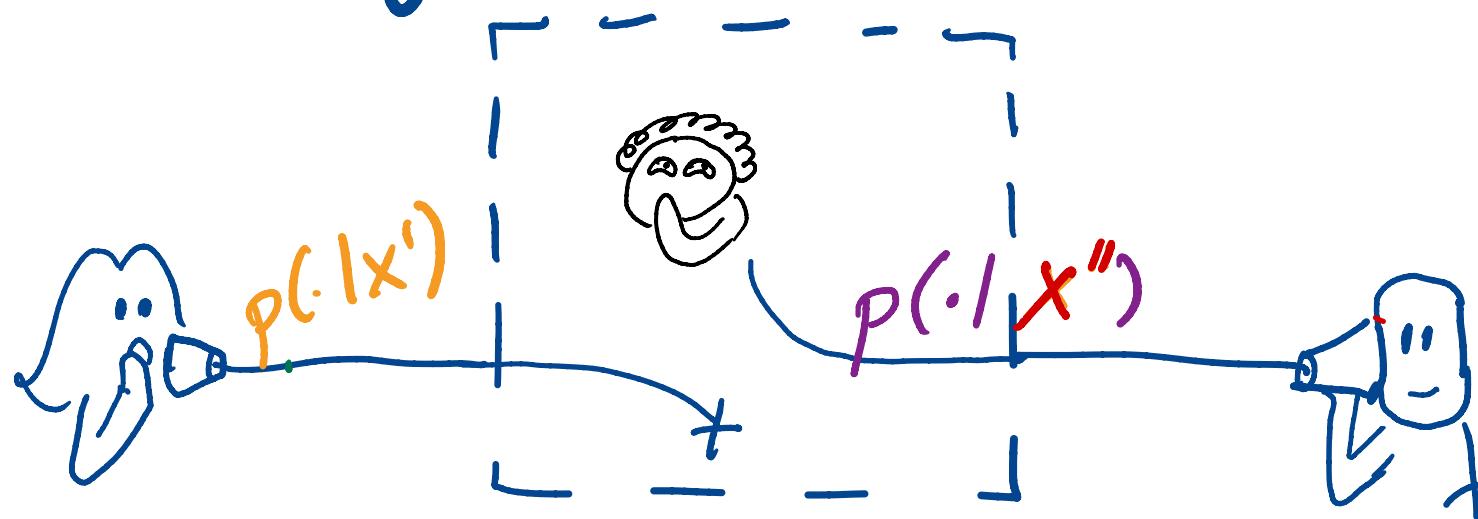


$p(\cdot|x_3)$ $p(\cdot|x_4)$ $p(\cdot|x')$
 $p(\cdot|x_2)$ $p(\cdot|x_1)$

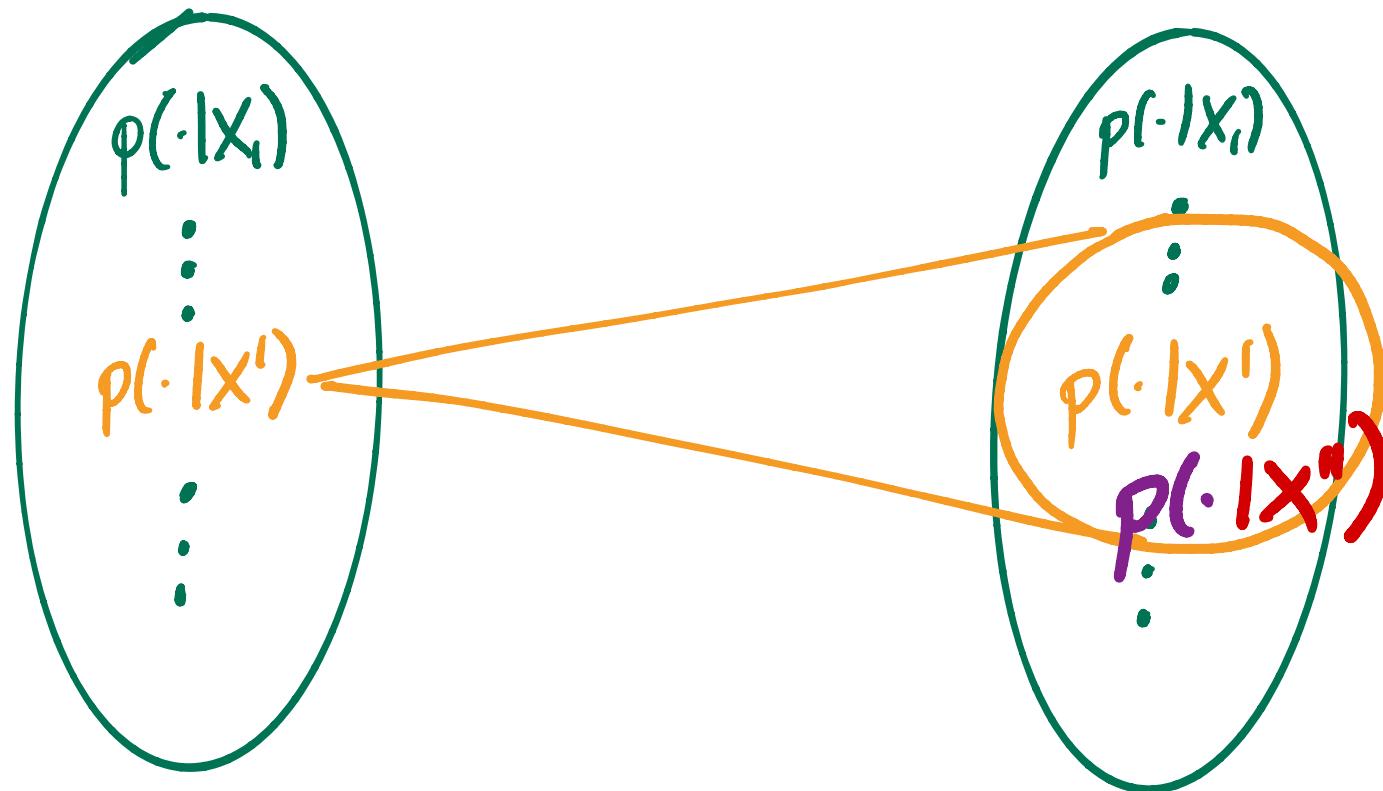
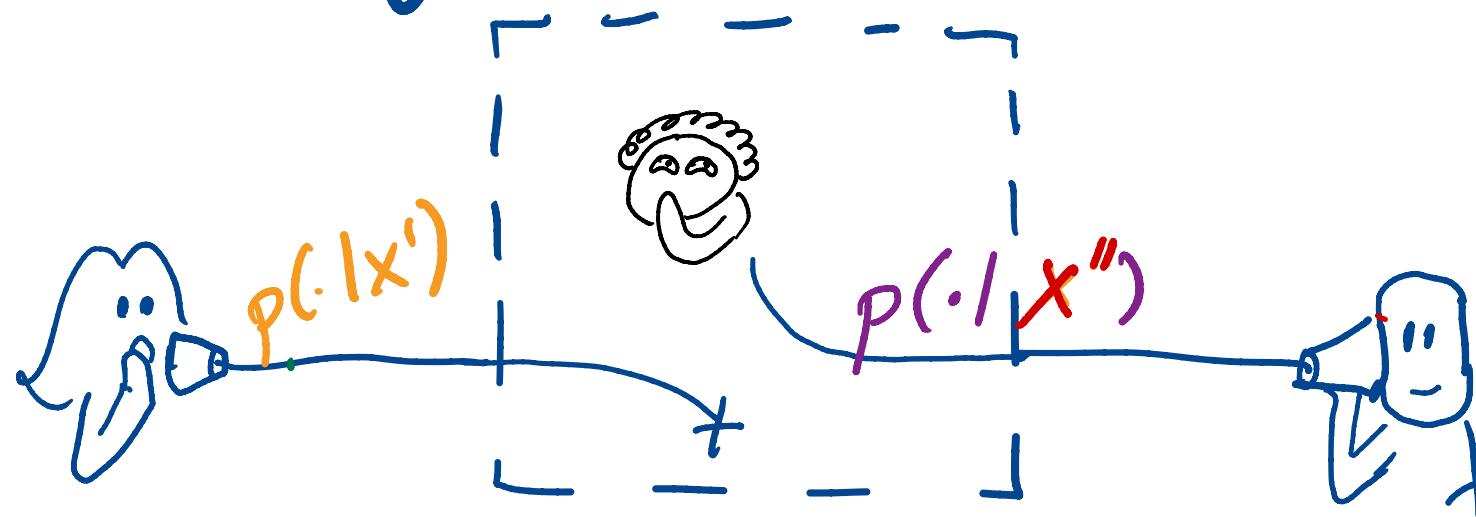
codeword space

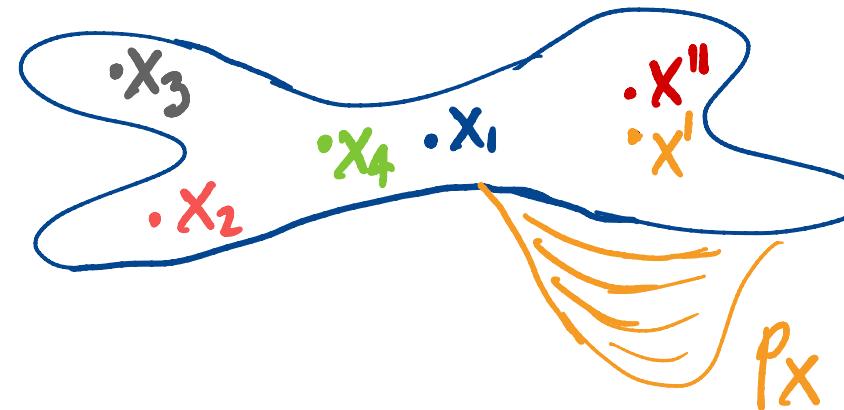
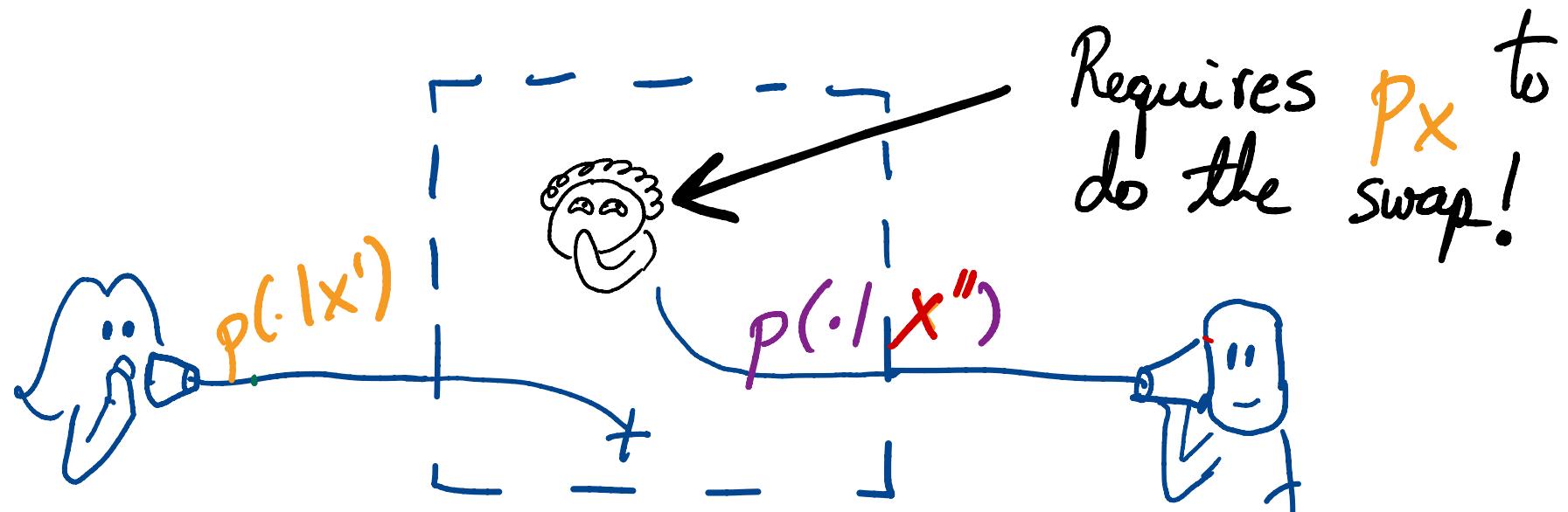
Instance space X

If the algo is robust...

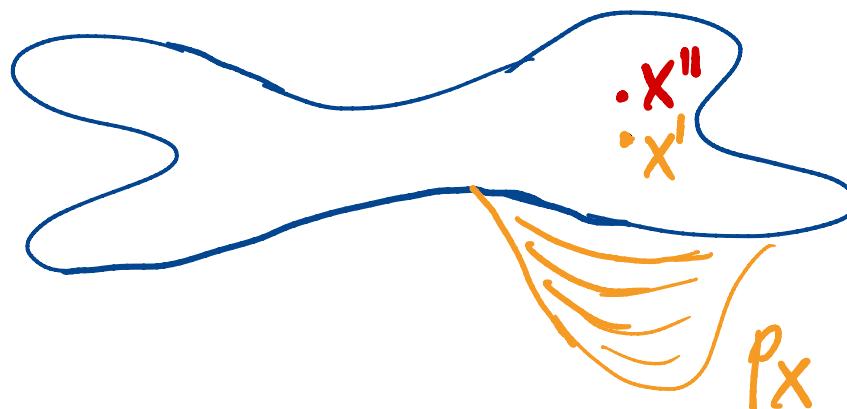
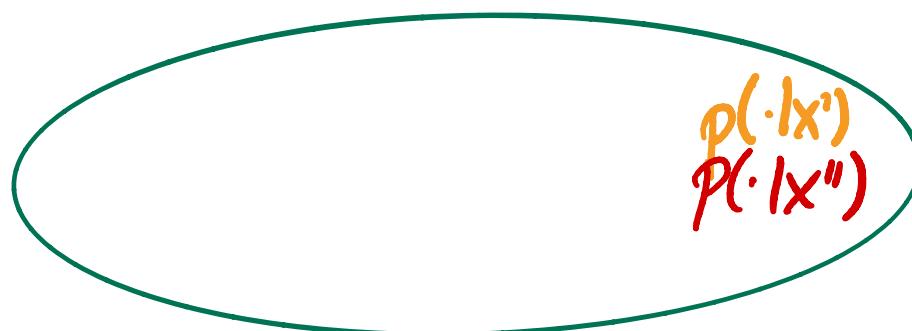
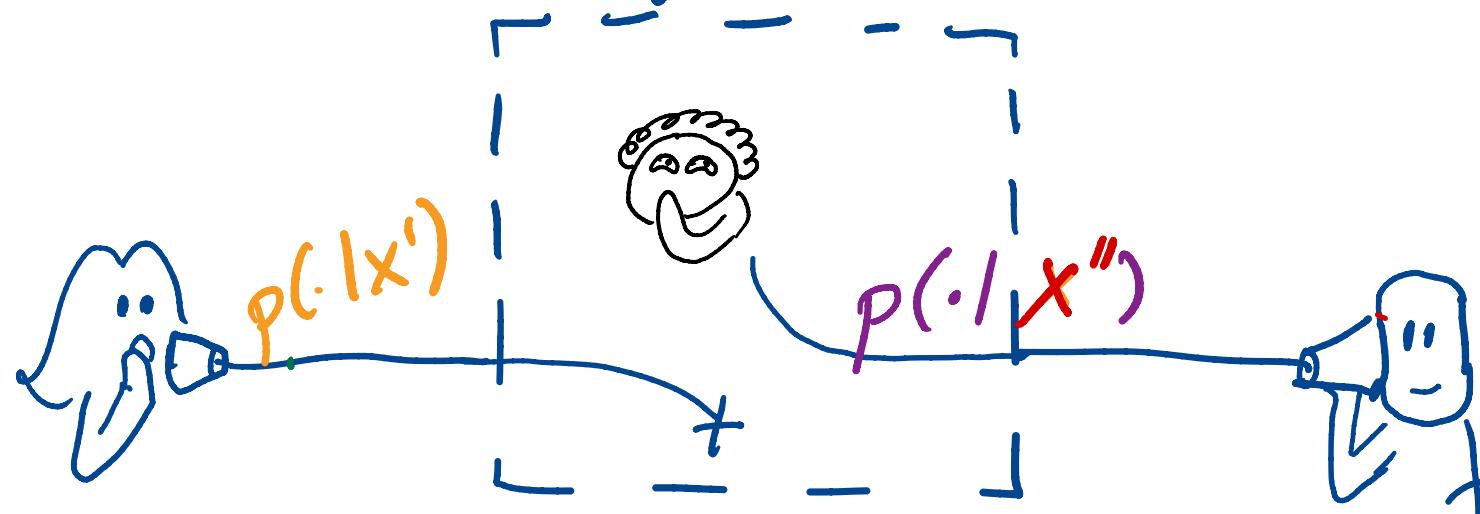


If the algo is *not* robust...



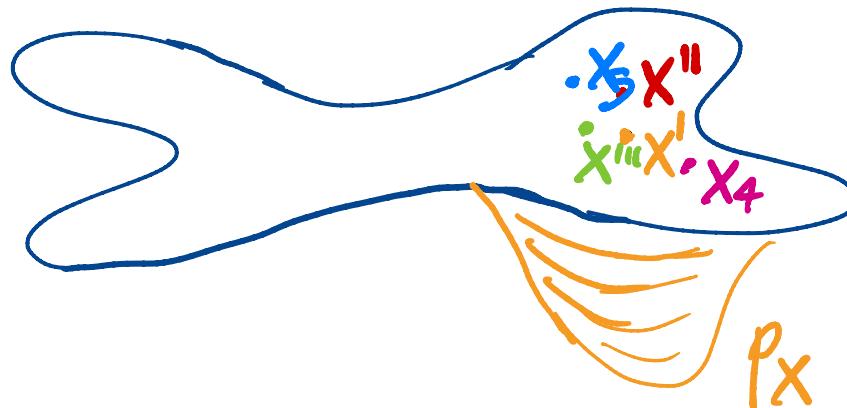
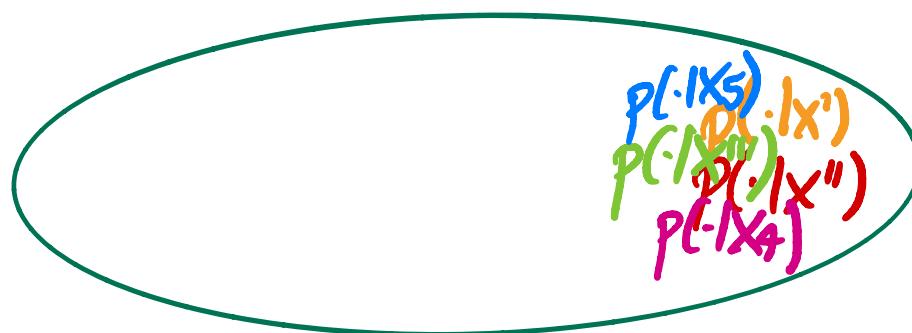
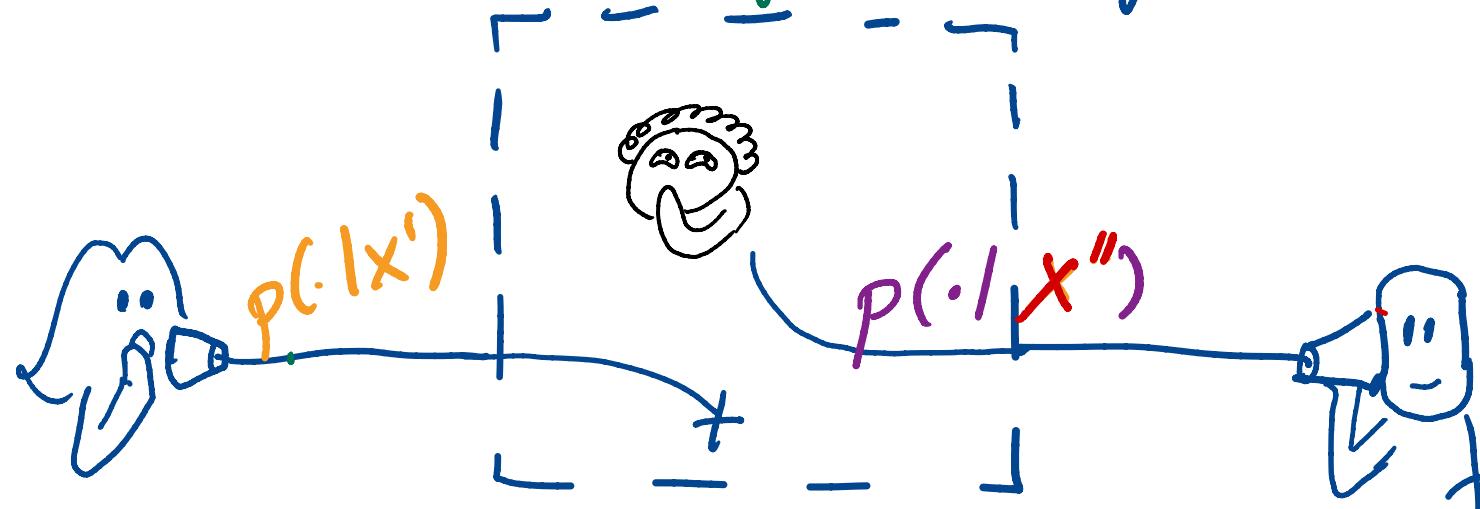


Problem: We only have one experiment P_X

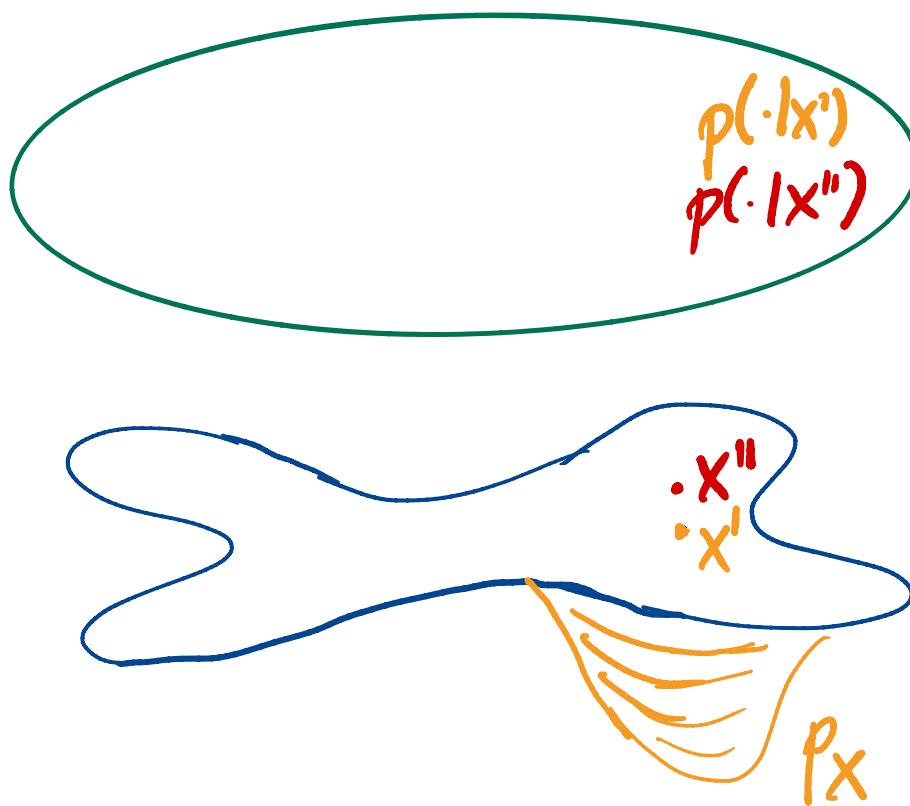


Problem:

Robust and bad-content algorithm }
Robust and good-content algorithm } yield similar
channels

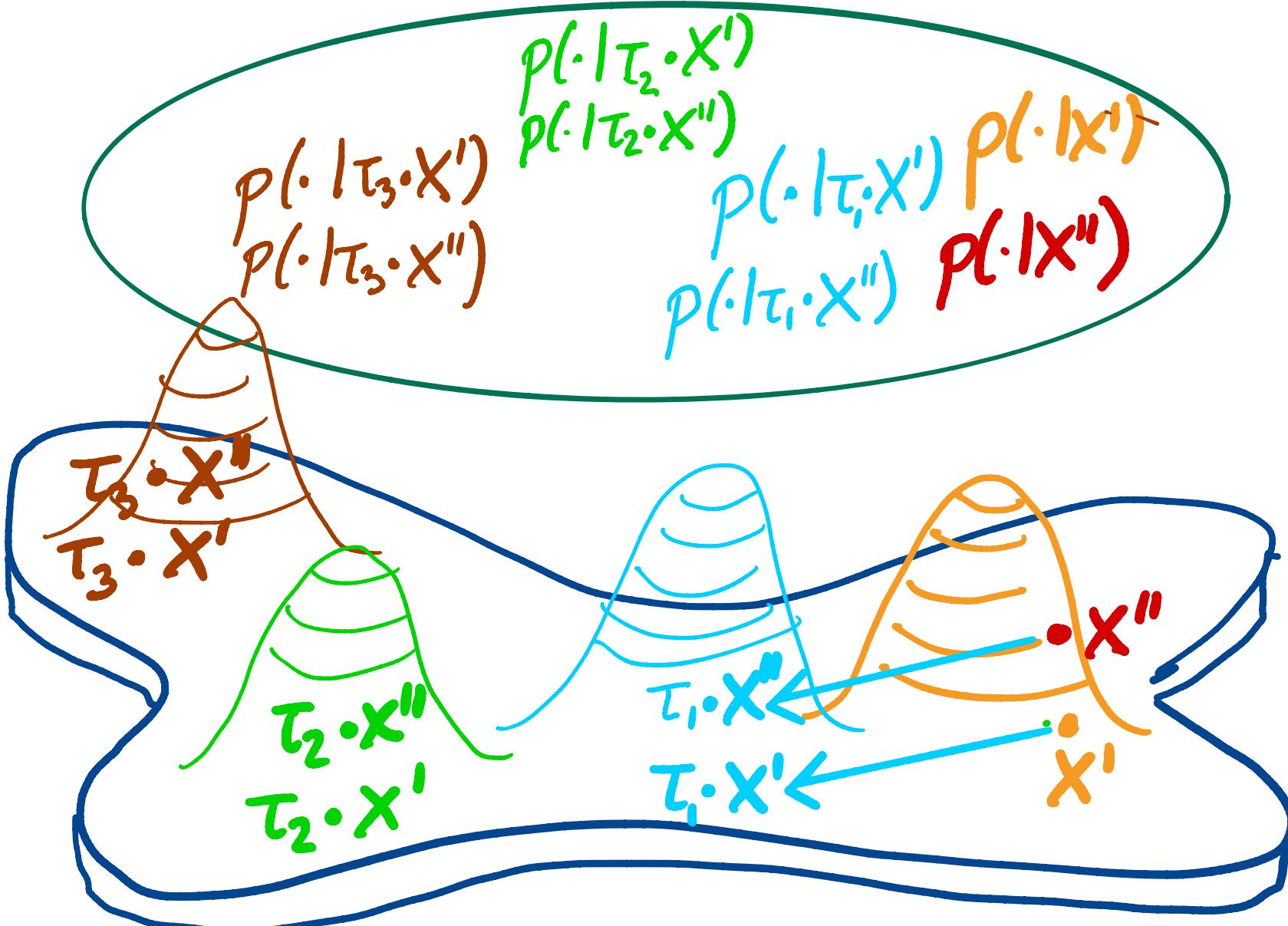


How to get more experiments ?



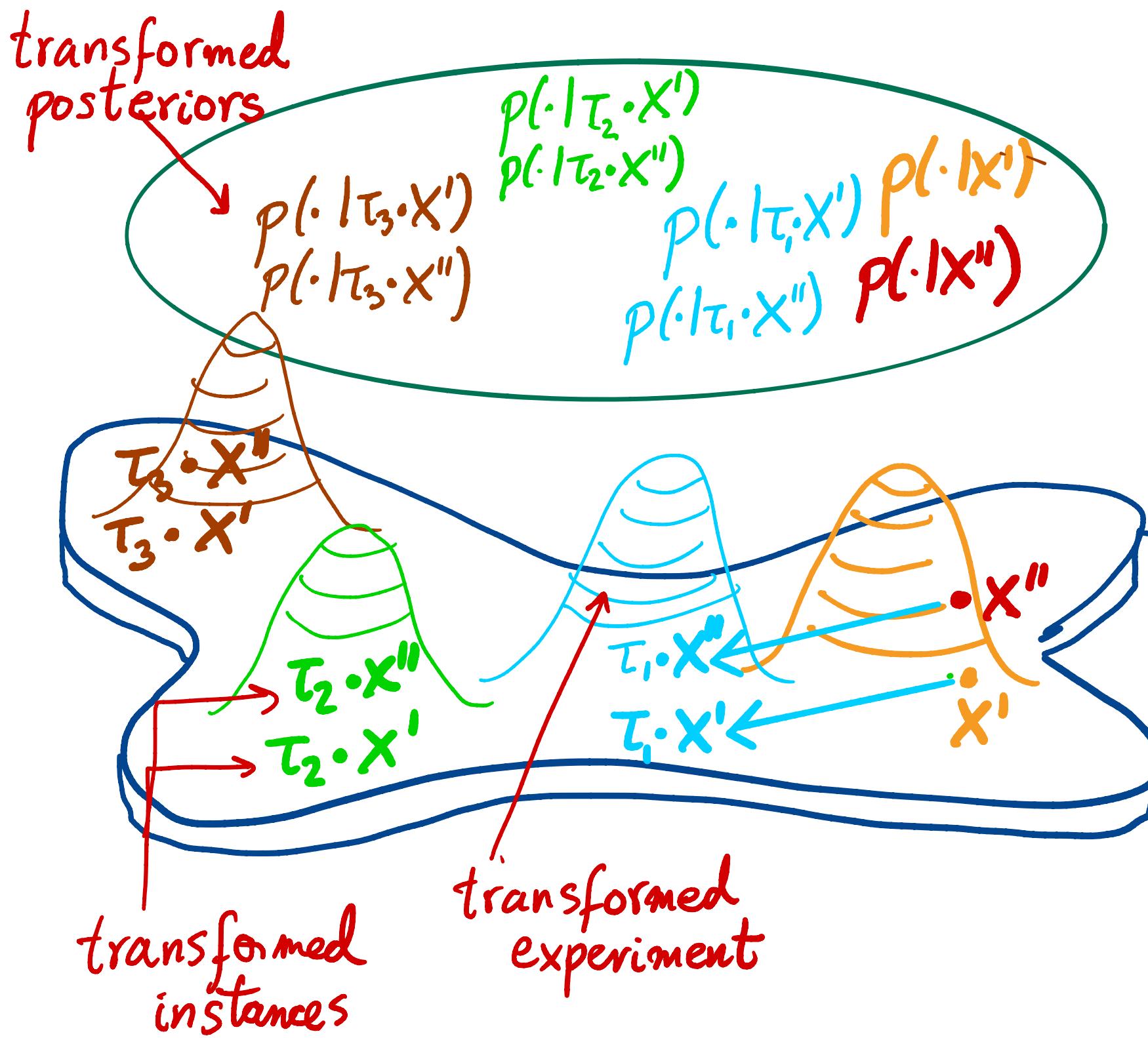
$p(\cdot|x')$
 $p(\cdot|x'')$

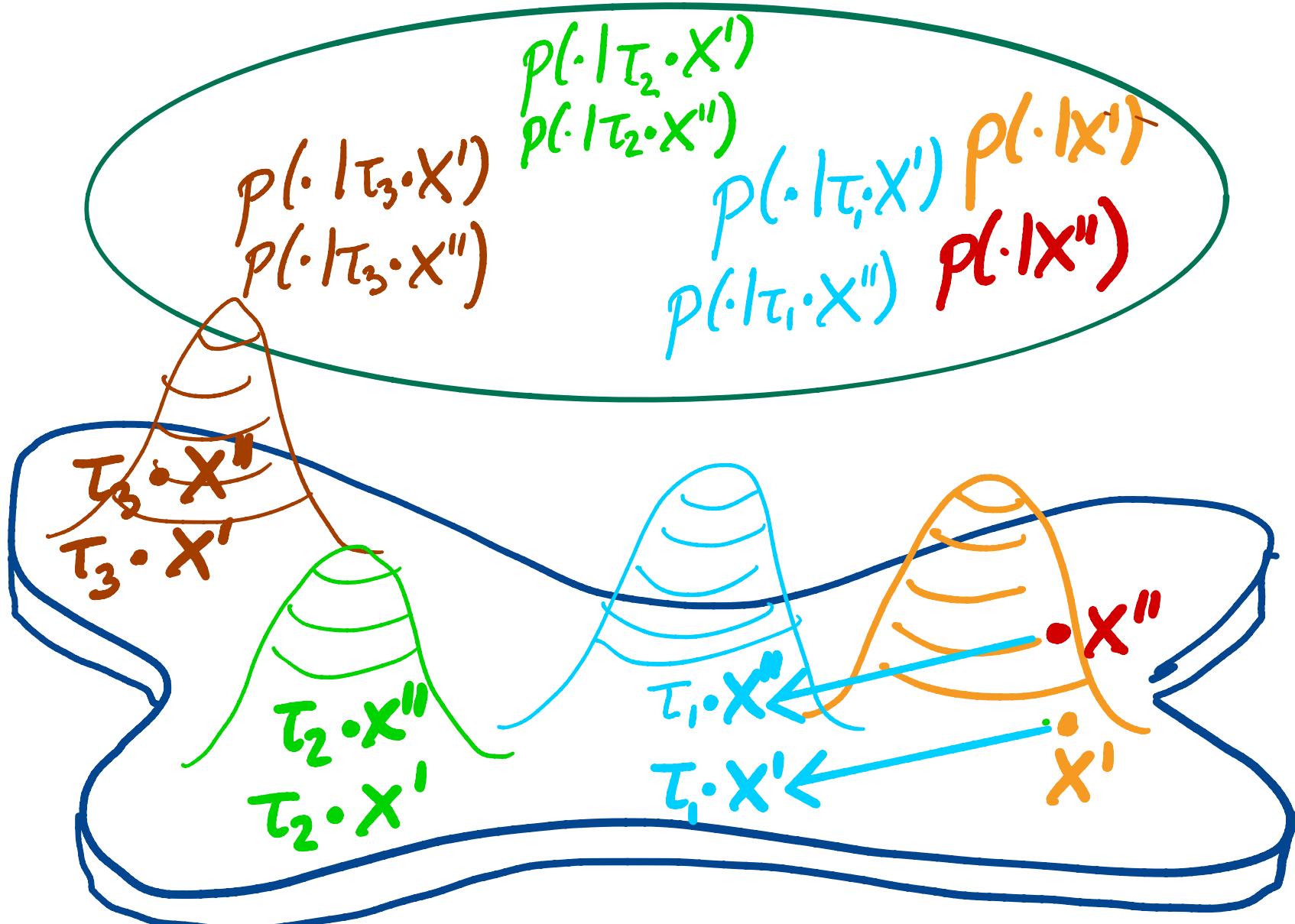




We define a series of transformations

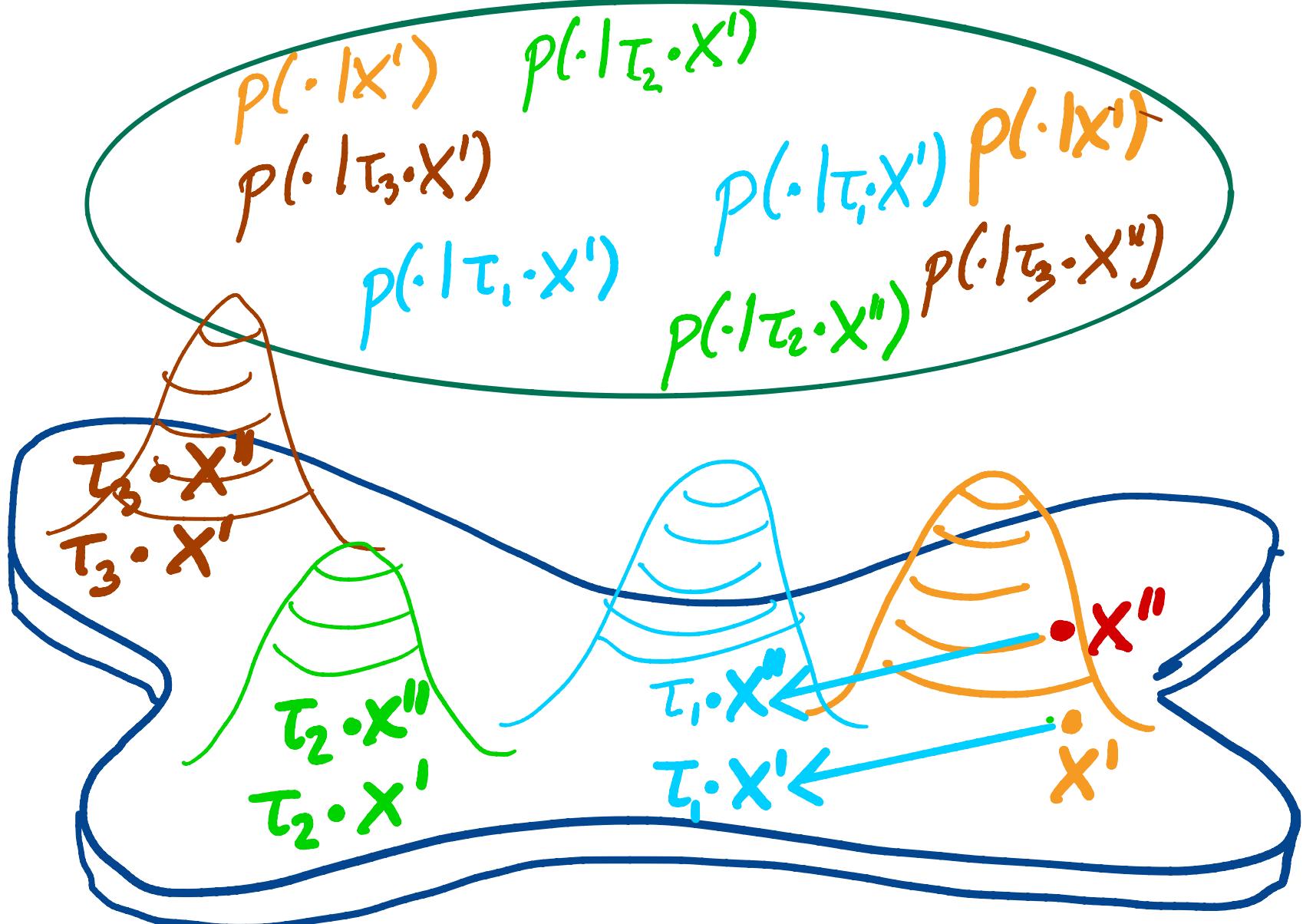
$$\mathcal{T} := \{T_1, T_2, T_3, \dots\}$$





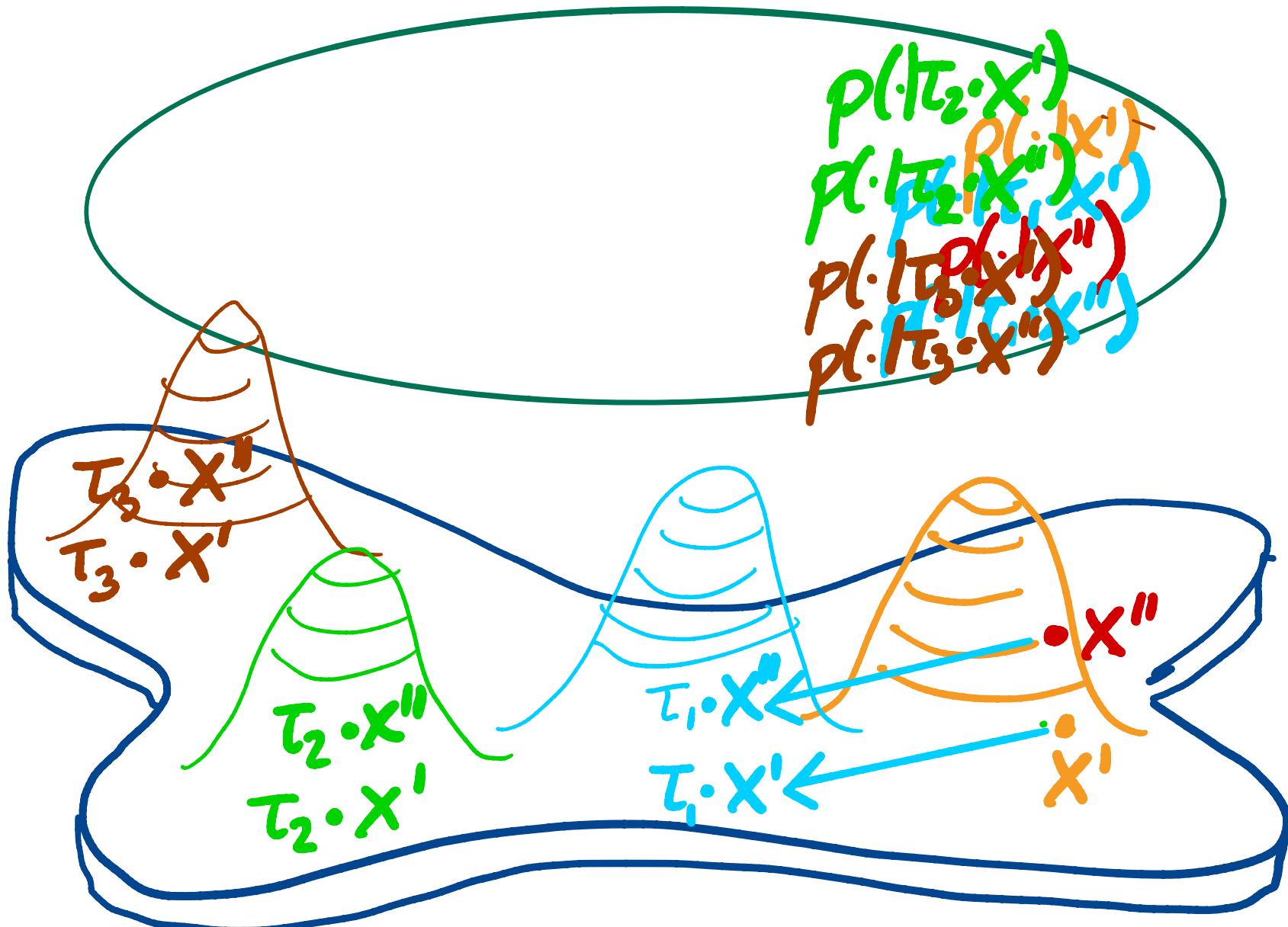
Robust and good-content algs:

- * different experiments \Rightarrow distant posteriors
- * same experiment \Rightarrow close posteriors



Unstable and good-content algs:

- * different experiments \Rightarrow distant posteriors
- * same experiment \Rightarrow distant posteriors



Robust and poor-content algs:

- * different experiments \Rightarrow close posteriors
- * same experiment \Rightarrow close posteriors

Requirements for transformations

The total mass of the transformed posteriors
should "uniformly" cover G .

That is $\sum_{\tau} p(c|\tau \cdot x^l) \approx \frac{|\Pi|}{|G|}$

Total mass of
 $p(\cdot|\tau_1 \cdot x^l), p(\cdot|\tau_2 \cdot x^l), \dots$

$$\sum_{\tau} \underbrace{\sum_c p(c|\tau \cdot x^l)}_1 = |\Pi|$$

$$* \sum_{\tau} p(c|\tau \cdot x^l) \in \left[\frac{|\Pi|}{|G|}(1-\rho), \frac{|\Pi|}{|G|}(1+\rho) \right],$$

for any $x^l \in X$

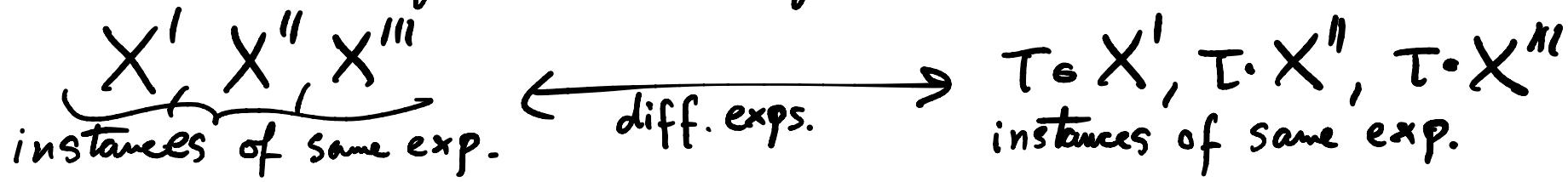
Requirements for transformations

* $\| p(\cdot | \tau_i \circ X^l) - p(\cdot | \tau_j \circ X^l) \|_1 > 0$, for $i \neq j$

That is, different transformations yield different posteriors

* The experimental measurements of X^l are not changed after a transformation.

That is, the noise in X^l is not changed after a transformation.



Examples of transformations

mean estimation

Translation δ $X \xrightarrow{\tau} X + \delta$

linear
regression

Scaling s $(X, Y) \xrightarrow{\tau} (X, sY)$

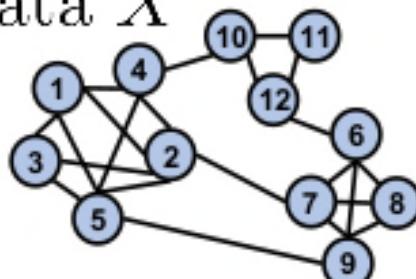
Graph-based
clustering

Vertex
permutation π

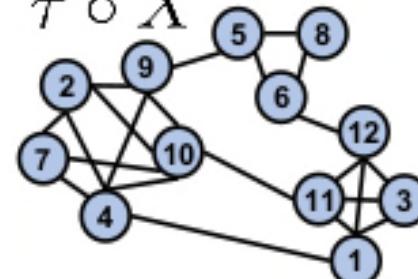
$(V, E) \xrightarrow{\tau} (V^\pi, E)$

$$V_i^\pi = V_{\pi(i)}$$

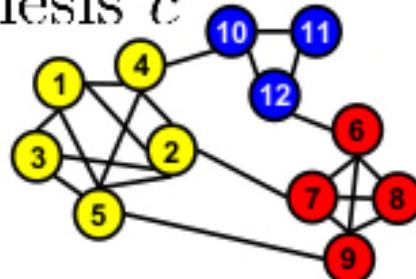
data X



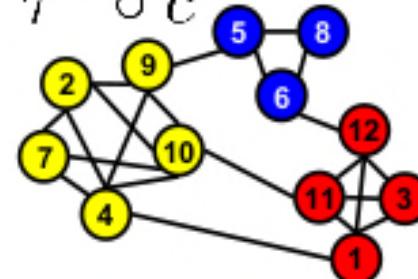
$\tau \circ X$



hypothesis c

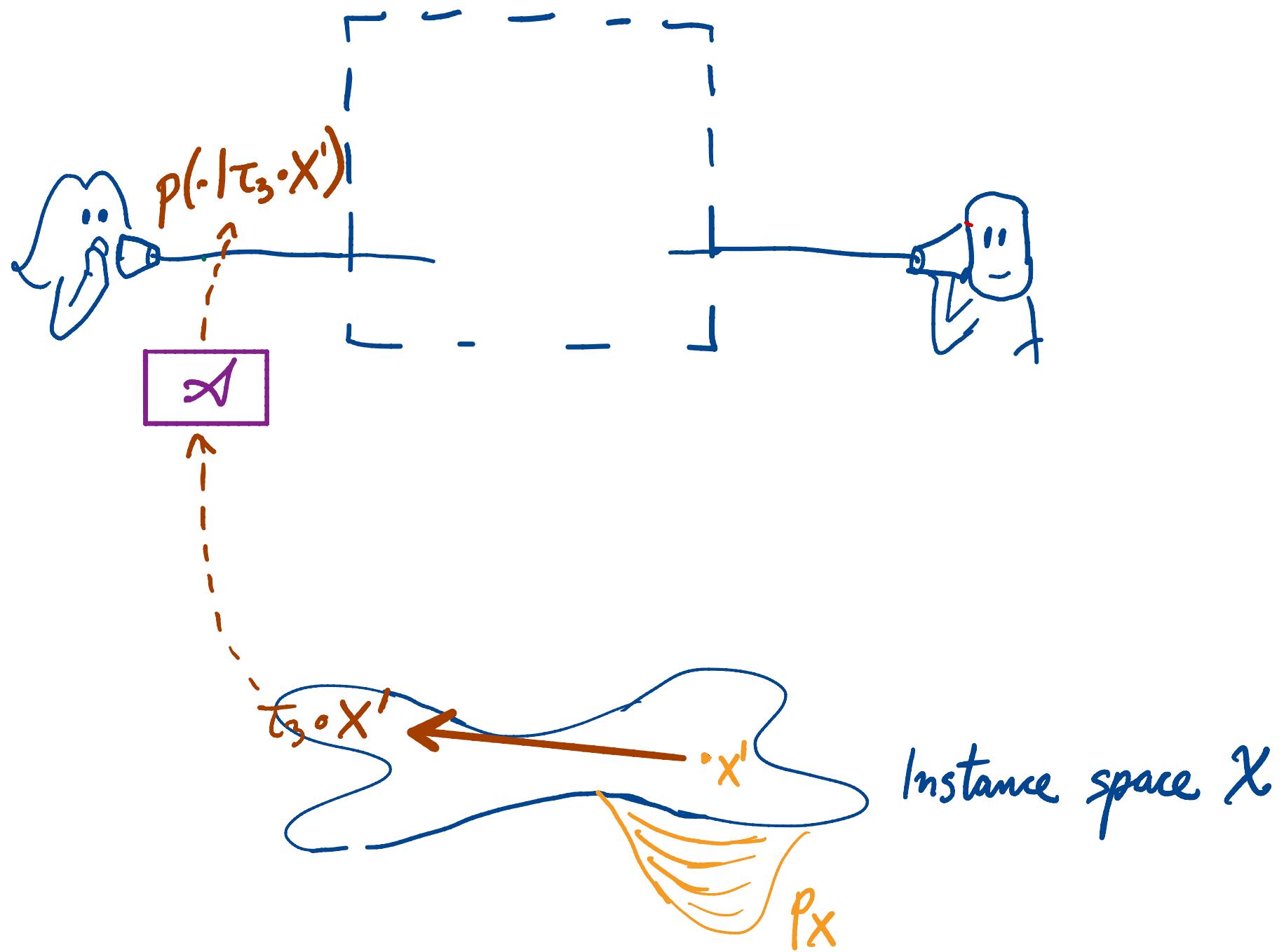


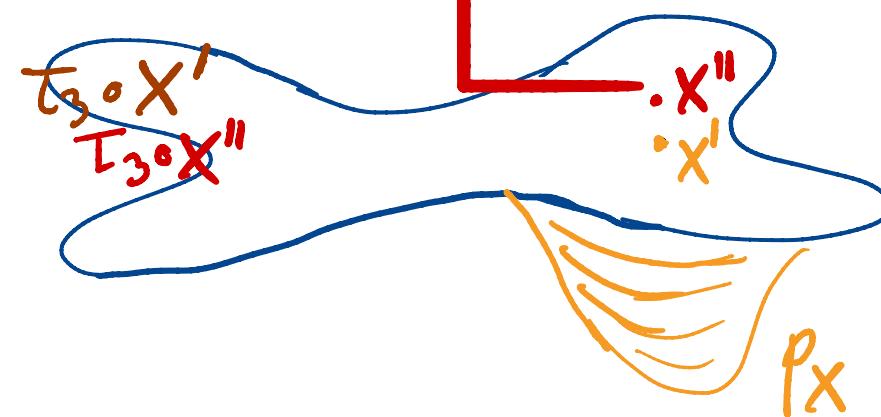
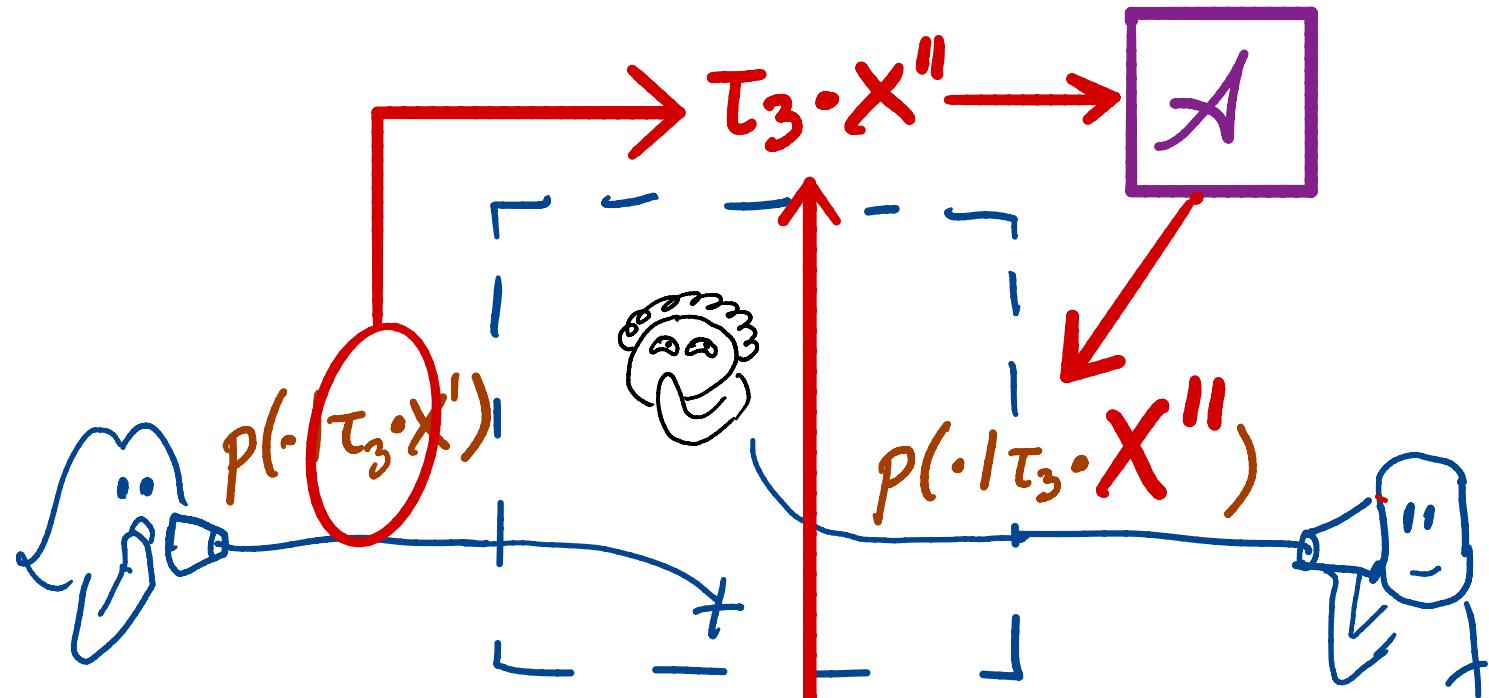
$\tau^c \circ c$



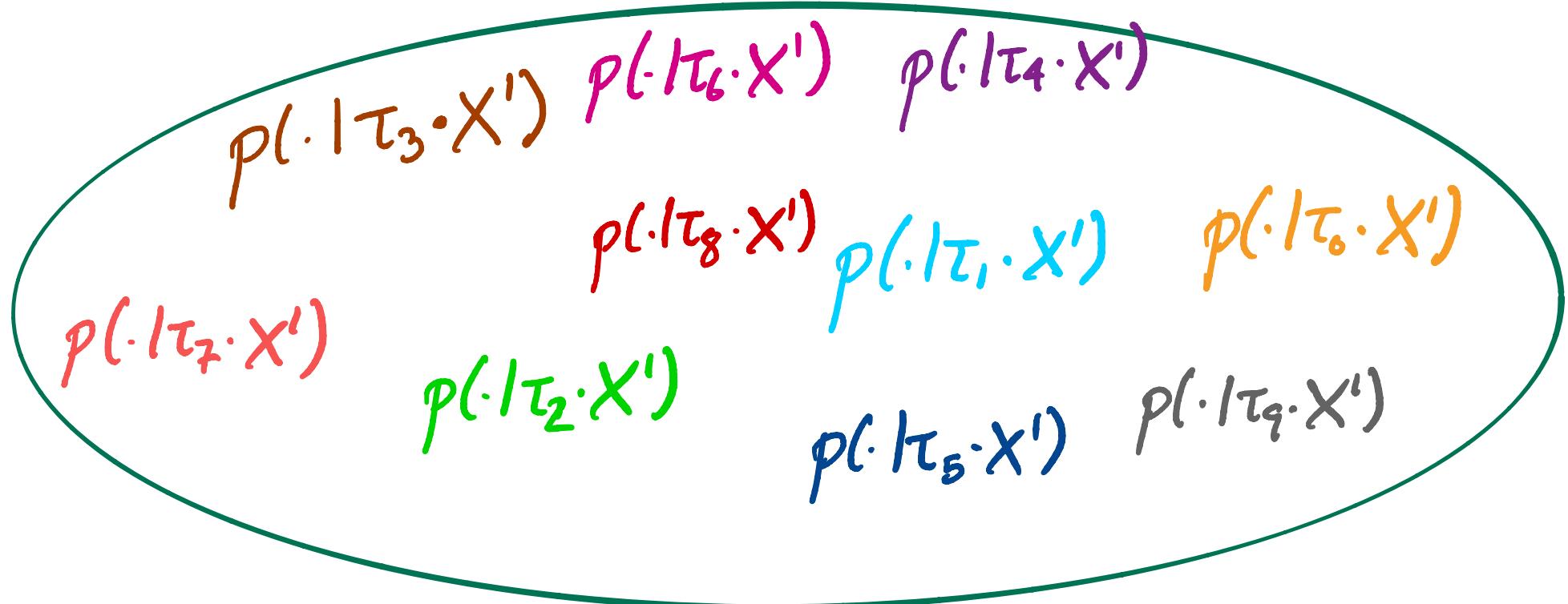
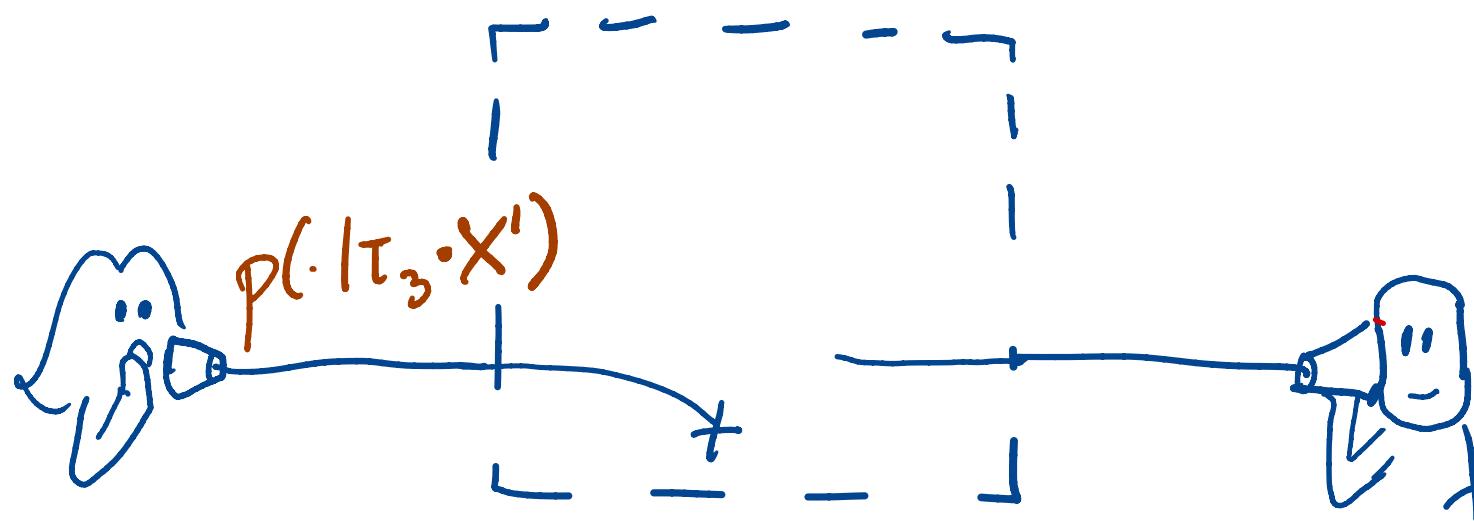
1 2 3 4 5 6 7 8 9 10 11 12

1 2 3 4 5 6 7 8 9 10 11 12

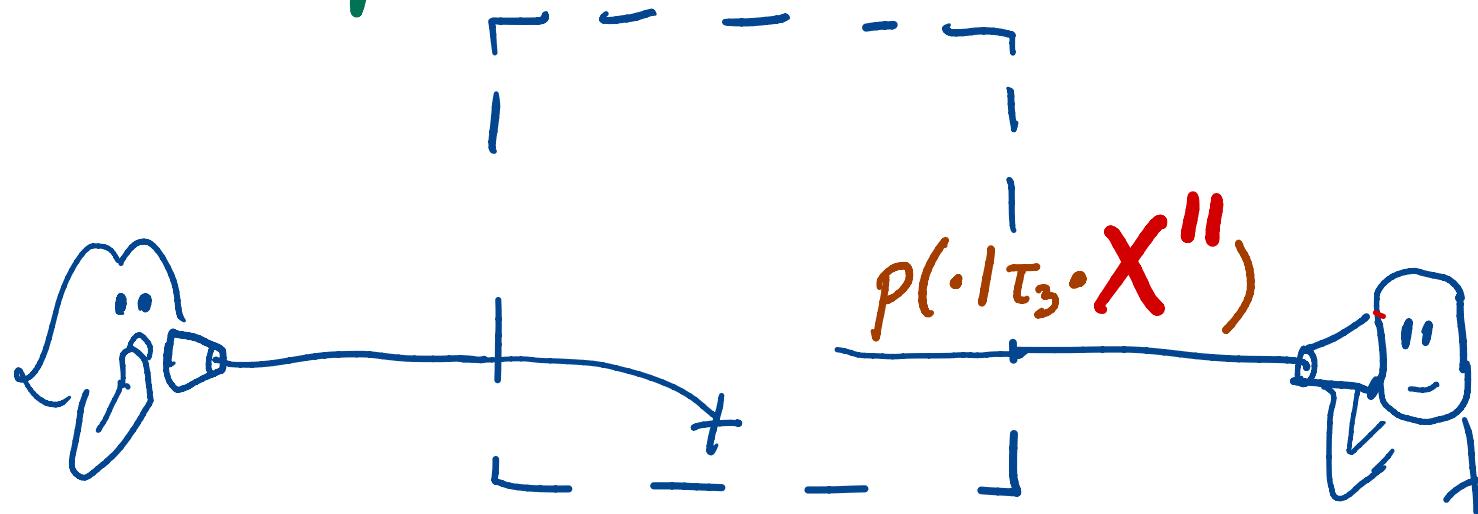




Input codeword space

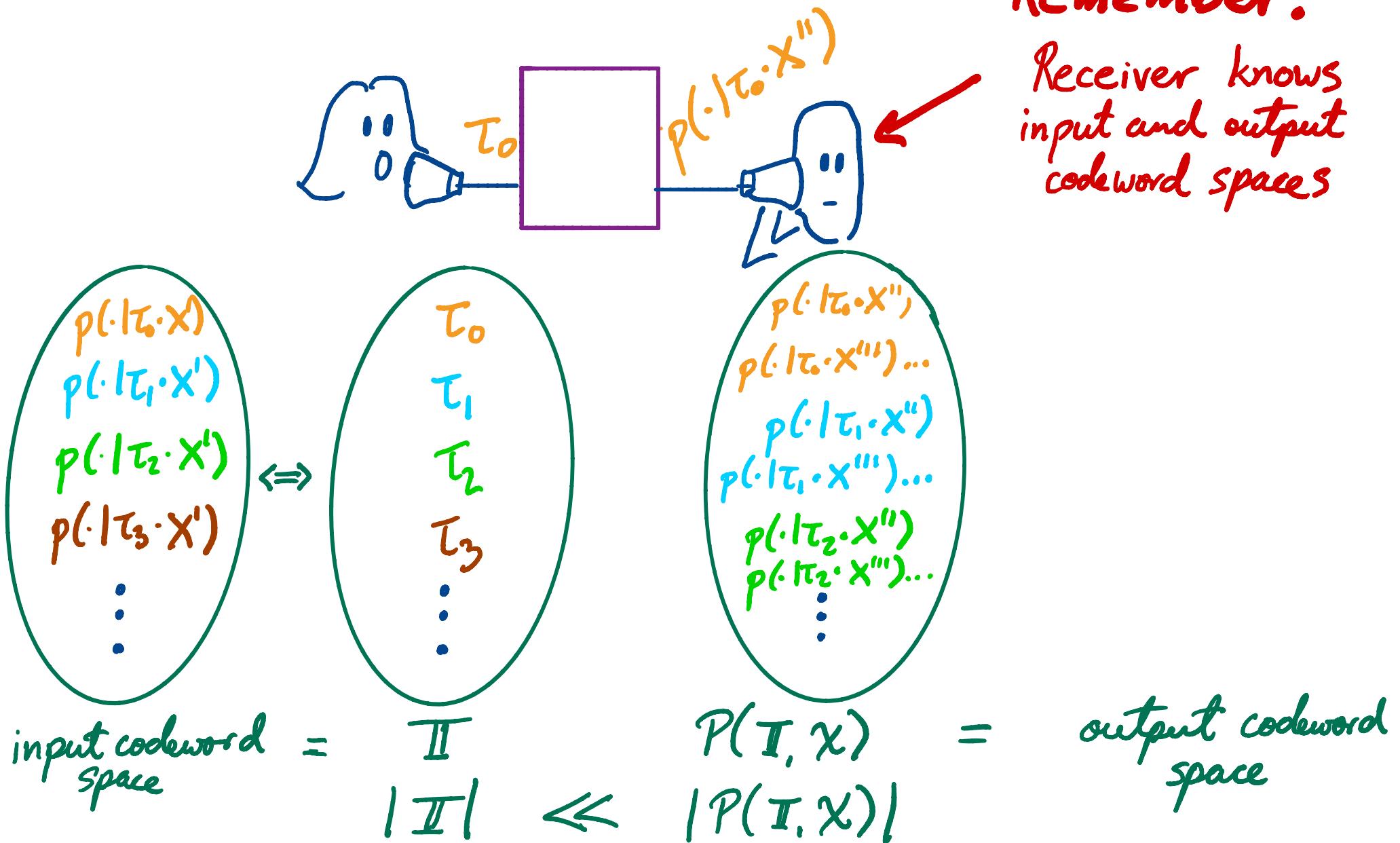


Output codeword space



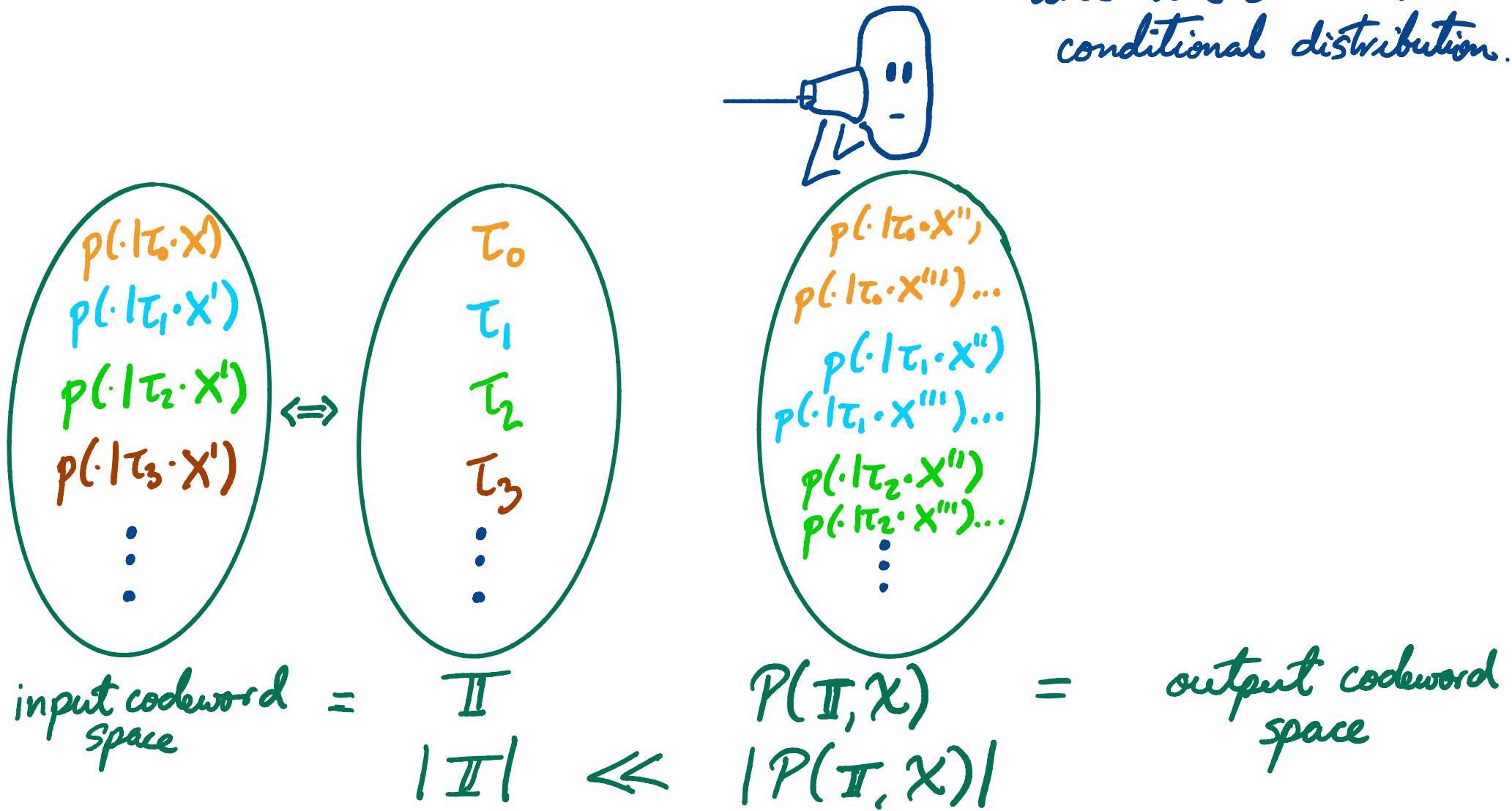
$p(\cdot | \tau_3 \cdot X'')$ $p(\cdot | \tau_4 \cdot X'')$
 $p(\cdot | \tau_7 \cdot X'')$ $p(\cdot | \tau_8 \cdot X'')$ $p(\cdot | \tau_0 \cdot X'')$
 $p(\cdot | \tau_2 \cdot X'')$ $p(\cdot | \tau_5 \cdot X'')$ $p(\cdot | \tau_9 \cdot X'')$
 $p(\cdot | \tau_1 \cdot X'')$

The algorithm at a channel



Receiver's knowledge

* Receiver knows X' , \mathcal{I} , A . That is, he knows input and output codeword spaces and the channel's conditional distribution.



The algorithm as a channel

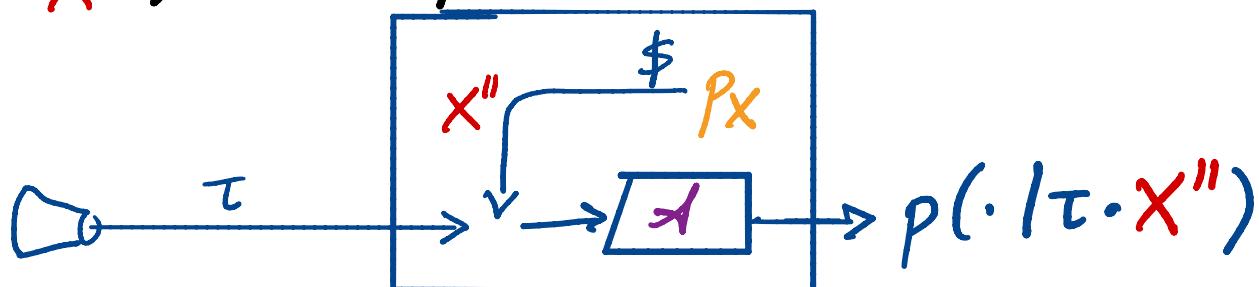
* Input alphabet :

$$\mathbb{I} = \{\tau_0, \tau_1, \dots, \tau_s\}.$$

* Output alphabet :

$$P(\mathbb{I}, \mathcal{X}) := \{p(\cdot | \tau \circ X'') : \tau \in \mathbb{I}, X'' \in \mathcal{X}\}$$

* When τ is input, a new instance X'' is drawn $\sim P_X$ and $p(\cdot | \tau \circ X'')$ is output.



Organization

What is PA?

Roadmap

Formalization of PA

→ Algorithms as channels

→ Derivation of PA

Applications

Derivation of PA

- * We follow Shannon's strategy.
Define a code for M messages by choosing M codewords at random.
- * We demonstrate that the prob. of a comm. error $\leq \text{const.} \cdot \exp(-\log |G| (\exp.\log. PA - \frac{\log M}{\log |G|} - \varepsilon))$
- * We compare this with Shannon's code's prob. of comm. error $\leq 2^{-n} (\text{capacity} - \frac{\log M}{n} - 3\varepsilon)$
- * Therefore, we propose the exp. log. PA as a capacity measure.

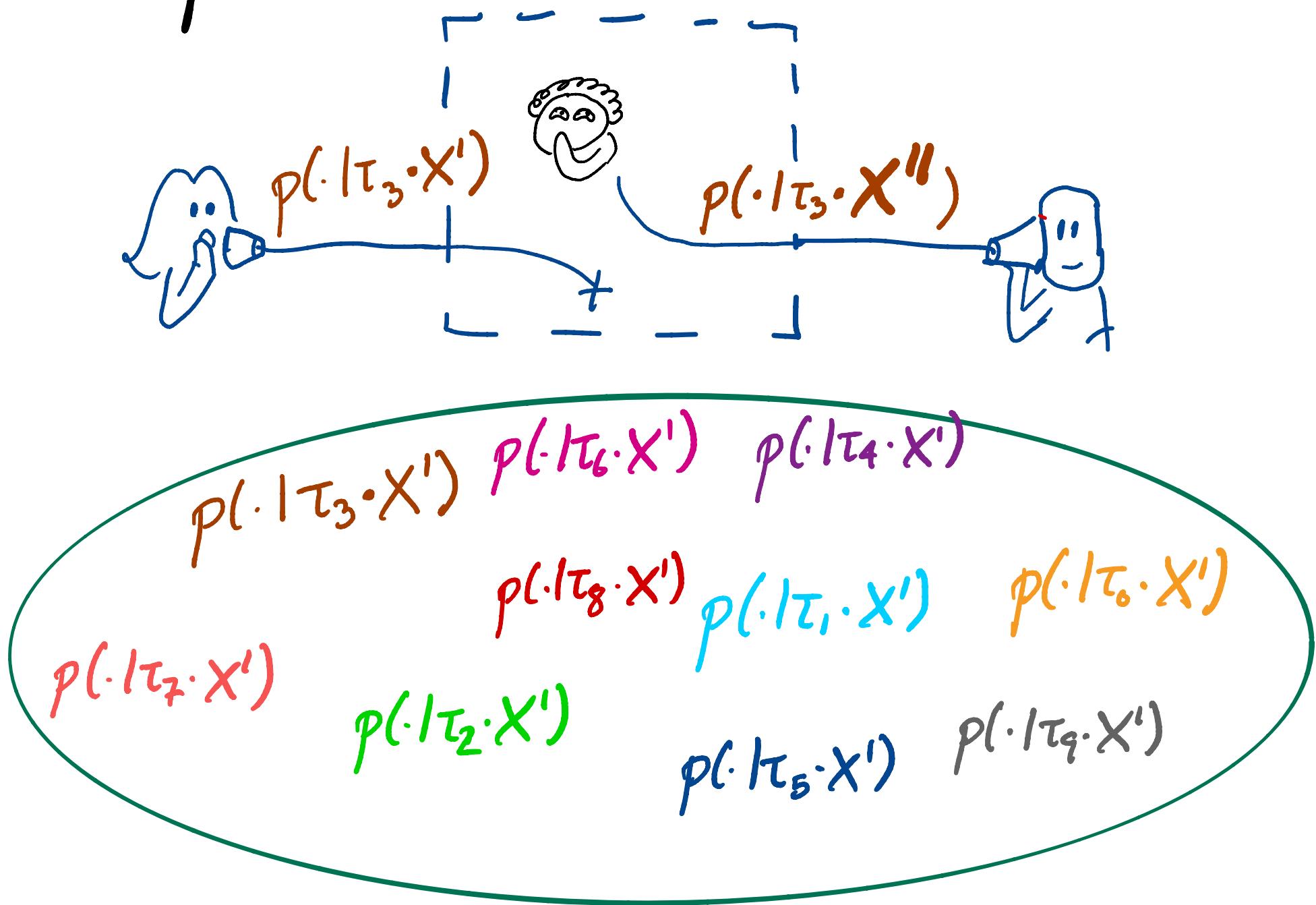
Remember:

A code is a pair (Enc, Dec) ,

$\text{Enc}: \{1 \dots M\} \longrightarrow$ ^{input} codeword space

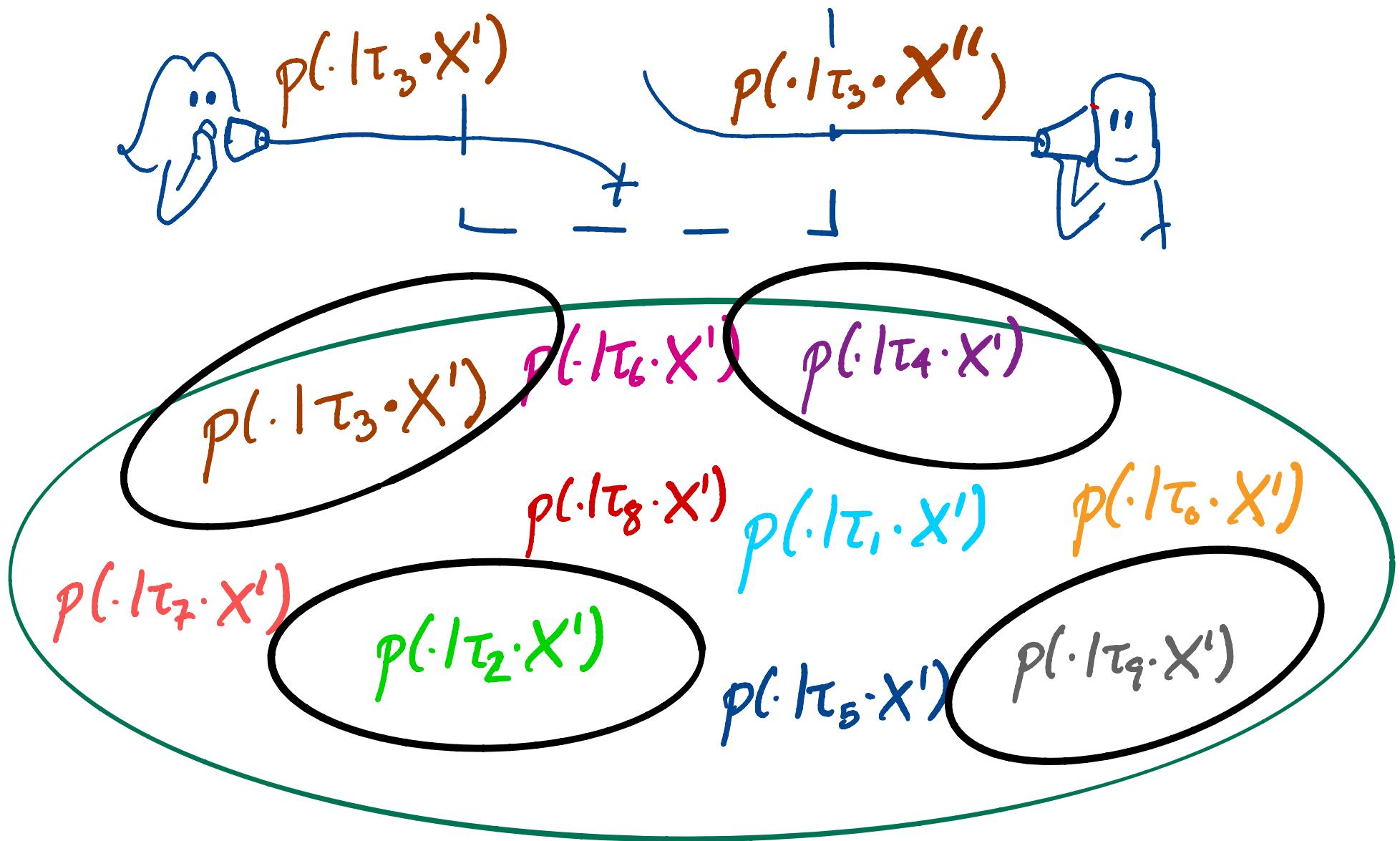
$\text{Dec}: \text{output codeword space} \longrightarrow \{1, \dots, M\}$

We follow Shannon's random code

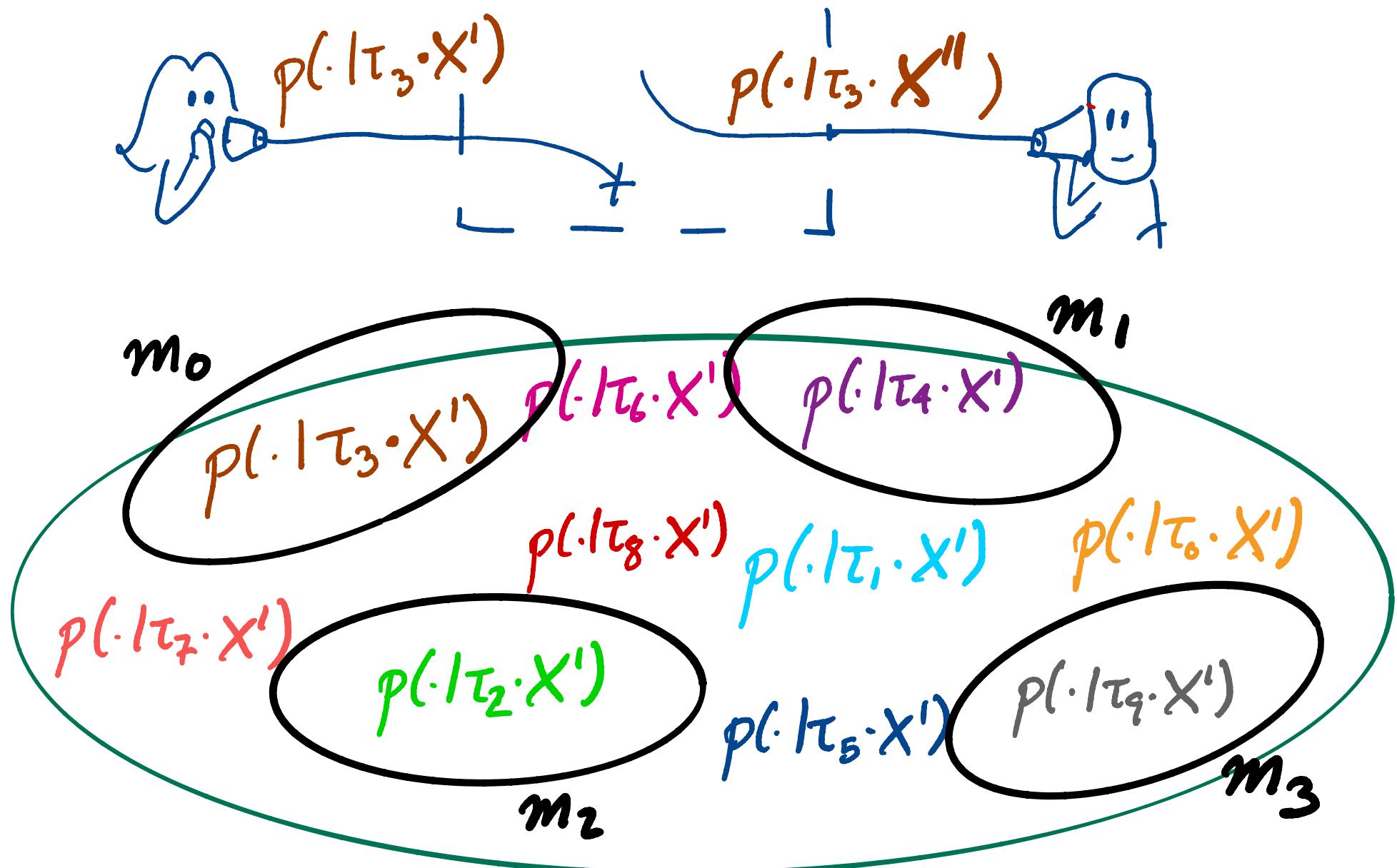


Let M be # msgs.

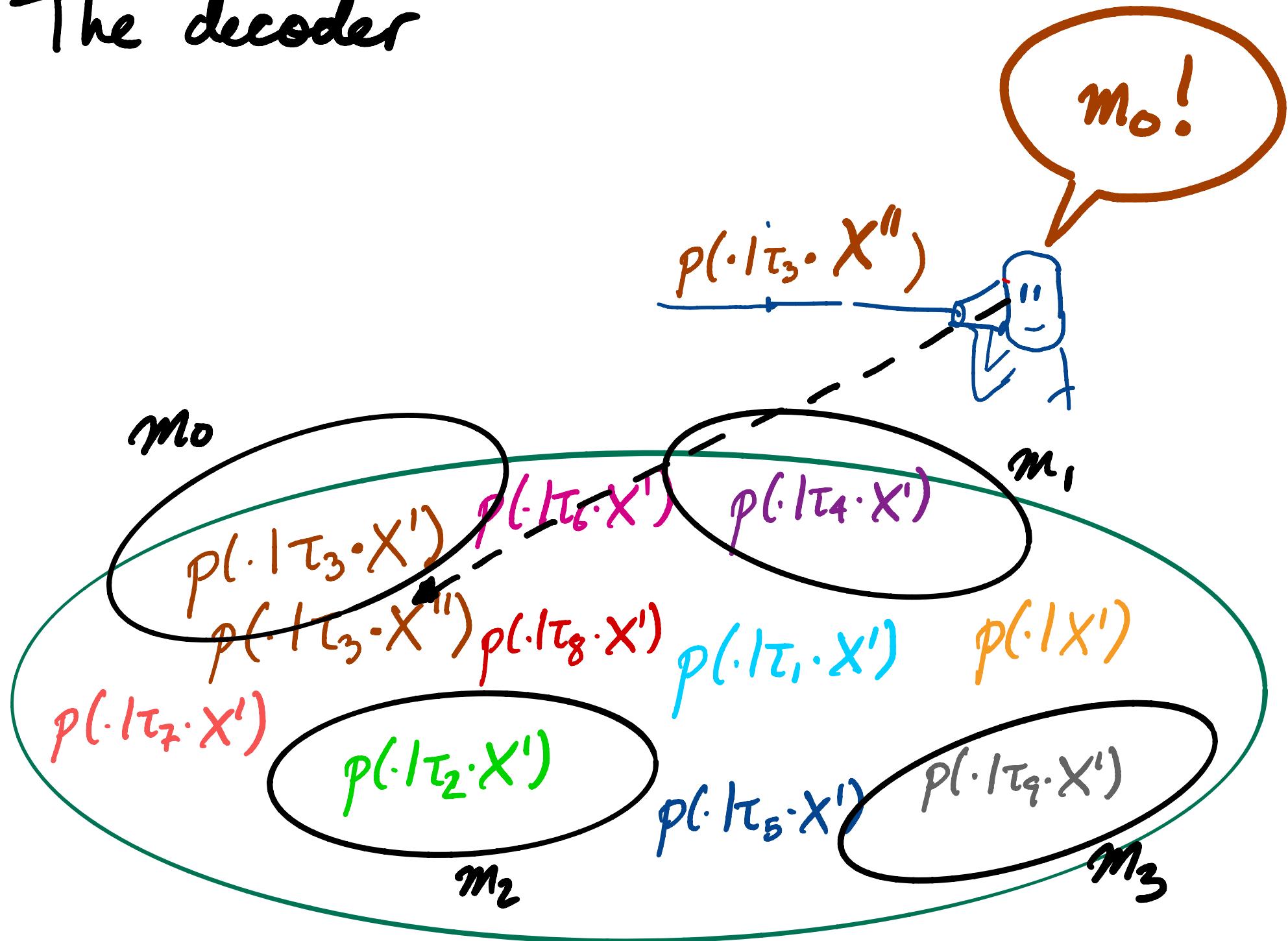
Pick M codewords at random.



The encoder function maps each message to each of these codewords.



The decoder

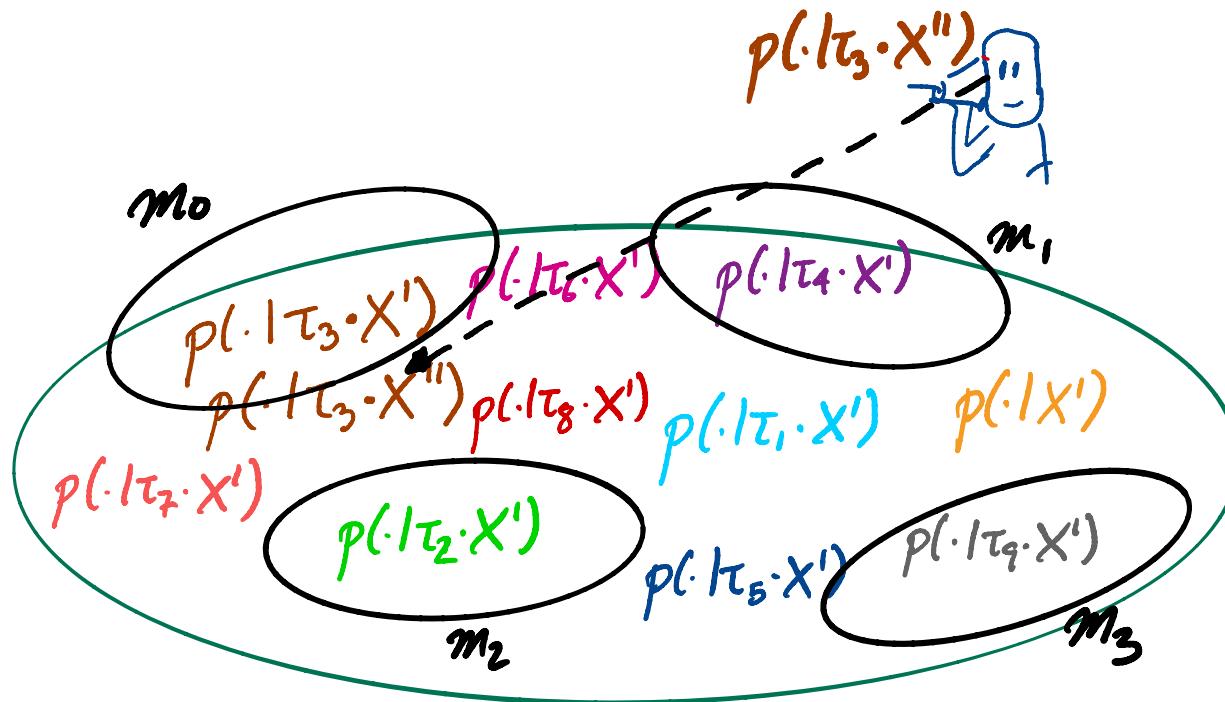


The decoder outputs

$$\arg \max_{m \leq M} \langle \text{Enc}(m), p(\cdot | \tau_3 \cdot X'') \rangle$$

$$= \arg \max_{\tau} \langle p(\cdot | \tau \cdot X'), p(\cdot | \tau_3 \cdot X'') \rangle$$

$$= \arg \max_{\tau} K(\tau \cdot X', \tau_3 \cdot X'')$$



Reference: Shannon's random code

- * Let $M \in \mathbb{N}$ and $X' \xleftarrow{\$} P_X$
- * Choose τ_1, \dots, τ_M from \mathcal{I} uniformly at random.
- * Sender and receiver learn $\{\tau_1, \dots, \tau_M\}$, X' , and \mathcal{A} .
- * Encoder $Enc: \{1 \dots M\} \rightarrow \{\tau_1 \dots \tau_m\}$.
 $i \longmapsto \tau_i$
- * Decoder $Dec_{X'}: P(\mathcal{I}, X) \rightarrow \{1, \dots, M\}$
 $p(\cdot | \tau \cdot X'') \longmapsto \hat{\tau}$

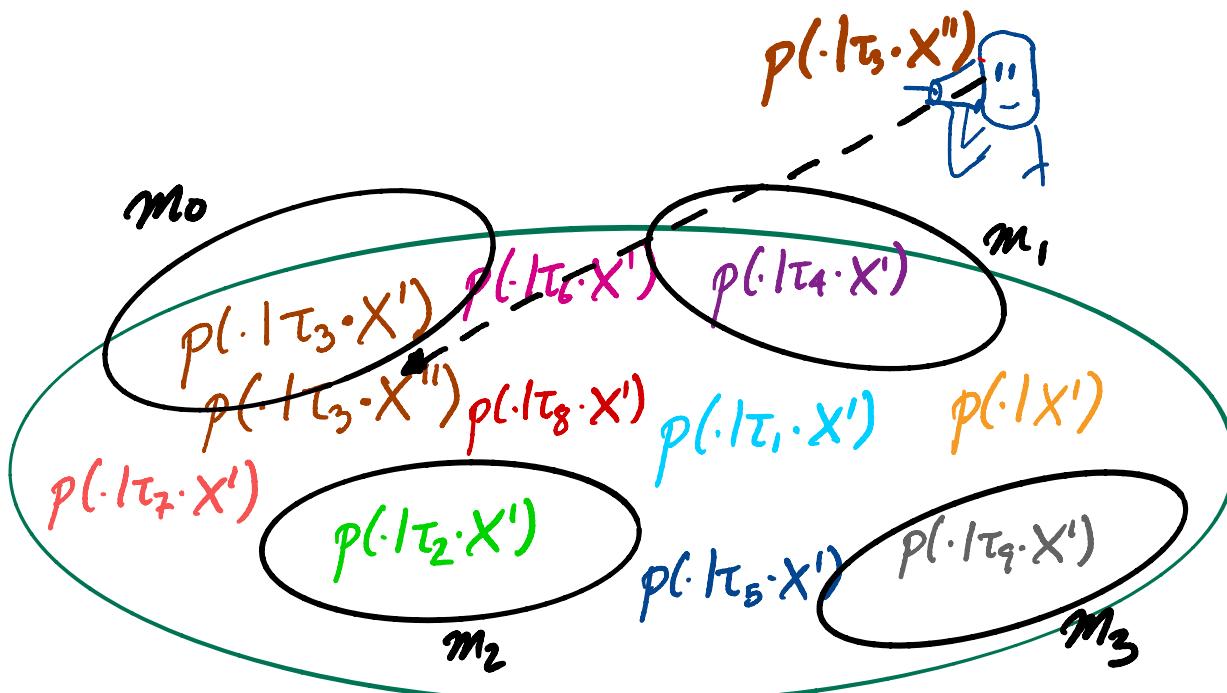
The receiver
needs X' , $\{\tau_1 \dots \tau_M\}$,
and \mathcal{A} to do this!

where $\hat{\tau} = \arg \max_{\tau' \in \{\tau_1 \dots \tau_M\}} K(\tau' \cdot X', \tau \cdot X'')$

Communication error

$$\arg \max_{\tau} K(\tau \circ X', \tau_3 \circ X'') \neq \tau_3$$

i.e., $K(\tau \cdot X', \tau_3 \cdot X'') \geq K(\tau_3 \cdot X', \tau_3 \cdot X'')$
for some $\tau \neq \tau_3$



Probability of a communication error

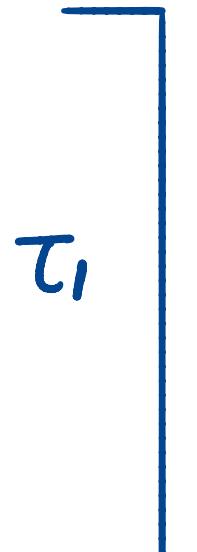
$$\Pr(E) = \Pr(E \mid \tau_1)$$

↑
Sender sent $\tau_1 = \text{id}$

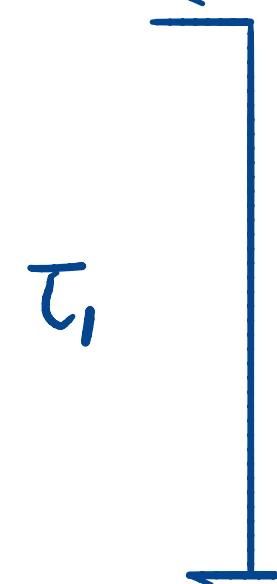
$$= \Pr(\kappa(\hat{\tau} \circ X', \tau_1 \circ X'') \geq \kappa(\tau_1 \circ X', \tau_1 \circ X''), \hat{\tau} \neq \tau_1 \mid \tau_1)$$
$$= \Pr \left[\begin{array}{c} \kappa(\tau_2 \circ X', \tau_1 \circ X'') \geq \kappa(\tau_1 \circ X', \tau_1 \circ X'') \\ \vdots \\ \kappa(\tau_M \circ X', \tau_1 \circ X'') \geq \kappa(\tau_1 \circ X', \tau_1 \circ X'') \end{array} \mid \tau_1 \right]$$

Apply the union bound

$$= \Pr \left[\begin{array}{l} K(\tau_2 \circ X', \tau_1 \circ X'') \geq K(\tau_1 \circ X', \tau_1 \circ X'') \text{ or} \\ K(\tau_3 \circ X', \tau_1 \circ X'') \geq K(\tau_1 \circ X', \tau_1 \circ X'') \text{ or} \\ \vdots \\ K(\tau_M \circ X', \tau_1 \circ X'') \geq K(\tau_1 \circ X', \tau_1 \circ X'') \end{array} \right]$$



$$= \Pr \left[\begin{array}{l} K(\tau_2 \circ X', X'') \geq K(X', X'') \text{ or} \\ K(\tau_3 \circ X', X'') \geq K(X', X'') \text{ or} \\ \vdots \\ K(\tau_M \circ X', X'') \geq K(X', X'') \end{array} \right]$$



$$\leq \sum_{1 < j \leq M} \Pr(K(\tau_j \circ X', X'') \geq K(X', X'') | \tau_1)$$

Apply Markov's inequality

2 / 5

$$\leq \sum_{1 \leq j \leq M} \Pr_{\substack{\text{randomness: } X', X'', \tau_1, \dots, \tau_M}} [K(\tau_j \cdot X', X'') \geq K(X', X'') \mid \tau_1]$$

$$= \sum_{1 \leq j \leq M} \sum_{X', X''} p(X', X'') \Pr_{\substack{\text{randomness: } X', X'', \tau_1, \dots, \tau_M}} [K(\tau_j \cdot X', X'') \geq K(X', X'') \mid \tau_1, X', X''] \leq$$

$$\leq \sum_{1 \leq j \leq M} \sum_{X', X''} p(X', X'') \frac{\mathbb{E}_{\tau_j} [K(\tau_j \cdot X', X'') \mid \tau_1, X', X'']}{K(X', X'')}$$

$$= \sum_{1 \leq j \leq M} \mathbb{E}_{X', X''} \left[\frac{\mathbb{E}_{\tau_j} [K(\tau_j \cdot X', X'') \mid \tau_1, X', X'']}{K(X', X'')} \right]$$

Bound the numerator

3/5

$$\leq \sum_{1 < j \leq M} E_{X', X''} \left[\frac{E_{\tau_j} [\kappa(\tau_j \circ X', X'') | \tau_i, X', X'']}{\kappa(X', X'')} \right]$$

$$= E_{\tau_j} \left[\sum_c p(c | \tau_j \circ X') p(c | X'') | \tau_i, X', X'' \right]$$

$$= \sum_c p(c | X'') E_{\tau_j} [p(c | \tau_j \circ X') | X']$$

$$= \sum_c p(c | X'') \sum_{\tau} \frac{1}{|\Pi|} p(c | \tau \circ X') \quad \sum_{\tau} p(c | \tau \circ X') \leq \frac{|\Pi|}{|\mathcal{C}|} (1+p)$$

$$= \sum_c p(c | X'') \frac{1}{|\Pi|} \frac{|\Pi|}{|\mathcal{C}|} (1+p) = \frac{(1+p)}{|\mathcal{C}|} \sum_c p(c | X'') = \frac{(1+p)}{|\mathcal{C}|}$$

Do the sum and get $\log K(x', x'')$

4/5

$$\begin{aligned} &\leq \sum_{1 \leq j \leq M} \mathbb{E}_{x', x''} \left[\mathbb{E}_{\tau_j} [\kappa(\tau_j \circ x', x'') | \tau_i, x', x''] \right] \\ &\leq \sum_{1 \leq j \leq M} \mathbb{E}_{x', x''} \left[\frac{1 + \rho}{|C| \kappa(x', x'')} \right] \\ &= (1 + \rho) (N - 1) \mathbb{E}_{x', x''} \left[\frac{1}{|C| \kappa(x', x'')} \right] \\ &\leq (1 + \rho) \mathbb{E}_{x', x''} \left[\frac{M}{|C| \kappa(x', x'')} \right] \\ &= (1 + \rho) \mathbb{E}_{x', x''} [\exp(-\log(|C| \kappa(x', x''))) + \log M] \end{aligned}$$

Apply AEP

5/5

$$= (1+p) \mathbb{E}_{x', x''} [\exp(-\log(|C|K(x', x'')) + \log M)]$$

$$= (1+p) \mathbb{E}_{x', x''} [\exp(-\log |C| (\underbrace{\frac{1}{\log |C|} \log(|C|K(x', x'')) - \frac{\log M}{\log |C|})}_{\text{emp. log. PA}})]$$

$$= (1+p) \mathbb{E}_{x', x''} [\exp(-\log |C|) \left[\begin{array}{l} \text{emp. log. PA} - \text{exp. log. PA} \\ + \text{exp. log. PA} - \frac{\log M}{\log |C|} \end{array} \right]]$$

$$= (1+p) \mathbb{E}_{X', X''} \left[\exp \left[-\log |C| \left(\begin{array}{c} \text{emp. log. PA} - \text{exp. log. PA} \\ + \text{exp. log. PA} - \frac{\log M}{\log |C|} \end{array} \right) \right] \right]$$

Asymptotic equipartition property:

$$\text{emp. log. PA} \xrightarrow[n \rightarrow \infty]{\text{in prob.}} \text{exp. log. PA}$$

$$\begin{aligned} & -\log |C| \text{emp. log. PA} + \log |C| \text{exp. log. PA} \\ & \leq \log |C| |\text{emp. log. PA} - \text{exp. log. PA}| \leq \log |C| \varepsilon \end{aligned}$$

with prob. $1-\varepsilon$
for large n .

$$\begin{aligned} & \leq (1+p) \mathbb{E}_{X', X''} \left[\exp \left(-\log |C| \left(-\varepsilon + \text{exp. log. PA} - \frac{\log M}{\log |C|} \right) \right) \right] \\ & = (1+p) \exp \left(-\log |C| \left(\text{exp. log. PA} - \frac{\log M}{\log |C|} - \varepsilon \right) \right). \end{aligned}$$

with prob. $1-\varepsilon$
for large n .

$$\Pr(E) \leq \text{const} \cdot \exp(-\log |G| (\text{exp. log. PA} - \frac{\log M}{\log |G|} - \varepsilon))$$

if $\text{exp. log. PA} > \frac{\log M}{\log |G|}$, then

$$\Pr(E) \xrightarrow[n \rightarrow \infty]{\text{in prob.}} 0$$

codeword length

channel capacity

code rate

$$\Pr(E) \leq 2^{-n} \left(\text{capacity} - \frac{\log \# \text{msgs}}{n} - 3\varepsilon \right)$$

if $\text{cap} > \frac{\log 2^{\text{nr}}}{n} = r$, then

$$\Pr(E) \xrightarrow[n \rightarrow \infty]{} 0.$$

$$\Pr(E) \leq \text{const} \cdot \exp(-\log |G| (\text{exp. log. PA} - \frac{\log M}{\log |G|} - \varepsilon))$$

if $\text{exp. log. PA} > \frac{\log M}{\log |G|}$, then

$$\Pr(E) \xrightarrow[n \rightarrow \infty]{\text{in prob.}} 0$$

\Rightarrow Large exp. log. PA

\Rightarrow Large M

\Rightarrow more messages can be communicated

\Rightarrow the algorithm is more informative

Asymptotic equipartition property

emp. log. PA $\xrightarrow[n \rightarrow \infty]{\text{in prob.}}$ exp. log. PA. ?

→ Solutions can be decomposed into "components"

* cluster assign fun \longrightarrow assig for each point.

* graph algos \longrightarrow edge

e.g. MST
shortest paths
clique

→ Posterior can be factorized; one factor per component

possibly via MFA

Asymptotic equipartition property

→ Posteriors can be factorized ; one factor per component
 $O(\log |C|)$ components

$$\begin{aligned}\rightarrow \kappa(x', x'') &= \sum_c p(c|x') p(c|x'') \approx \sum_c \prod p_i(c_i|x') p_i(c_i|x'') \\ &\xrightarrow{\text{see lecture 4 (DA)}} = \prod_i \sum_g p_i(g|x') p_i(g|x'') \\ &= \prod_i \kappa_i(x', x'')\end{aligned}$$

Asymptotic equipartition property

$$\text{emp. log. PA} = \frac{1}{\log |C|} \log (|C| K(x', x''))$$

$$= 1 + \frac{1}{\log |C|} \log K(x', x'')$$

$$\approx 1 + \frac{1}{\log |C|} \log \prod_i K_i(x', x'')$$

$$= 1 + \frac{1}{\log |C|} \sum_i \log K_i(x', x'')$$

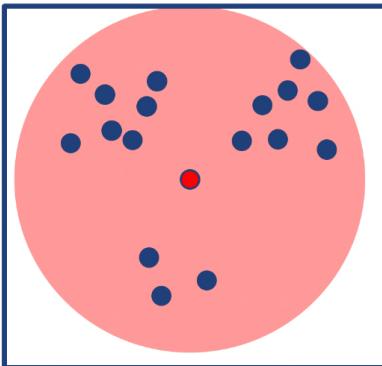
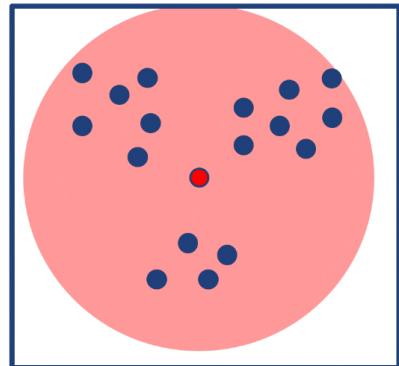
components
 $= O(\log |C|)$

in prob.
 $n \rightarrow \infty$

exp. log. PA.

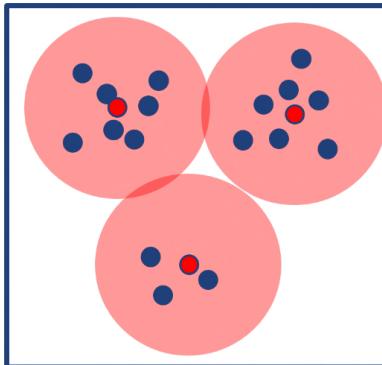
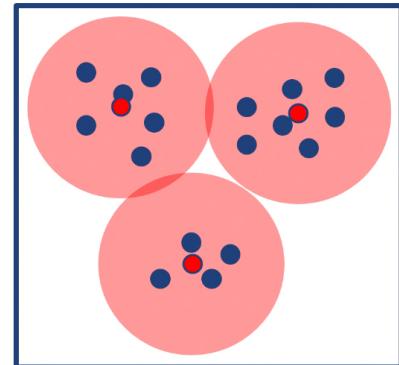
Summary

Algorithms must be robust and
retrieve signals



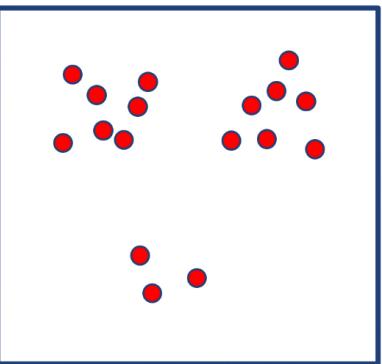
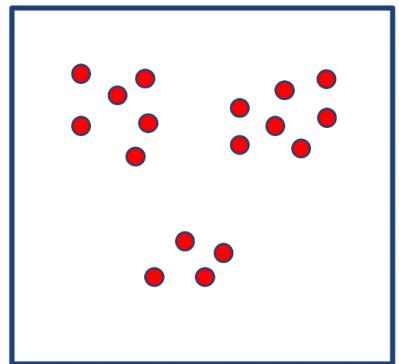
robust

no
signal



robust

signal



not
robust

signal

Information extraction

* Outputs have content

* Algorithm is robust

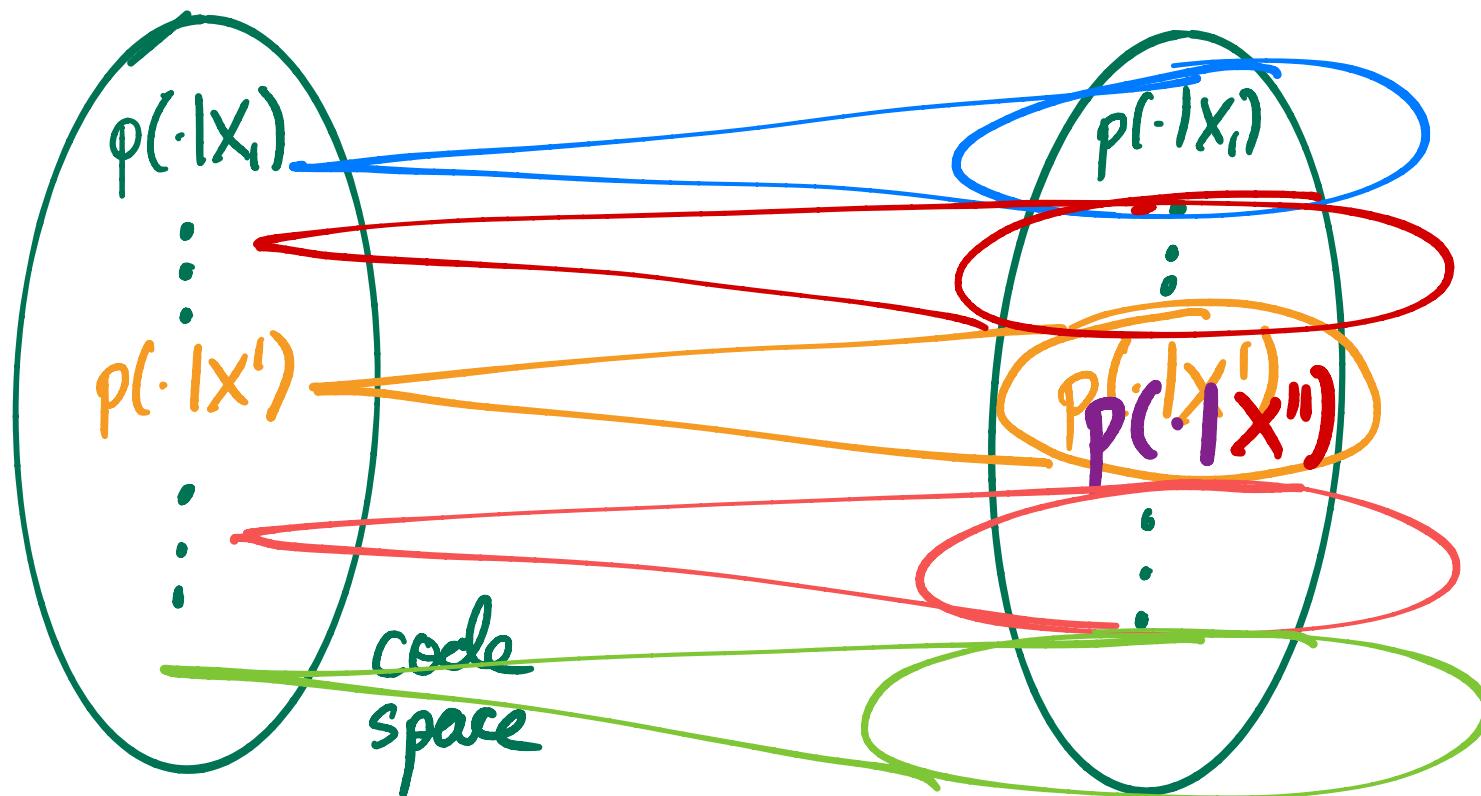
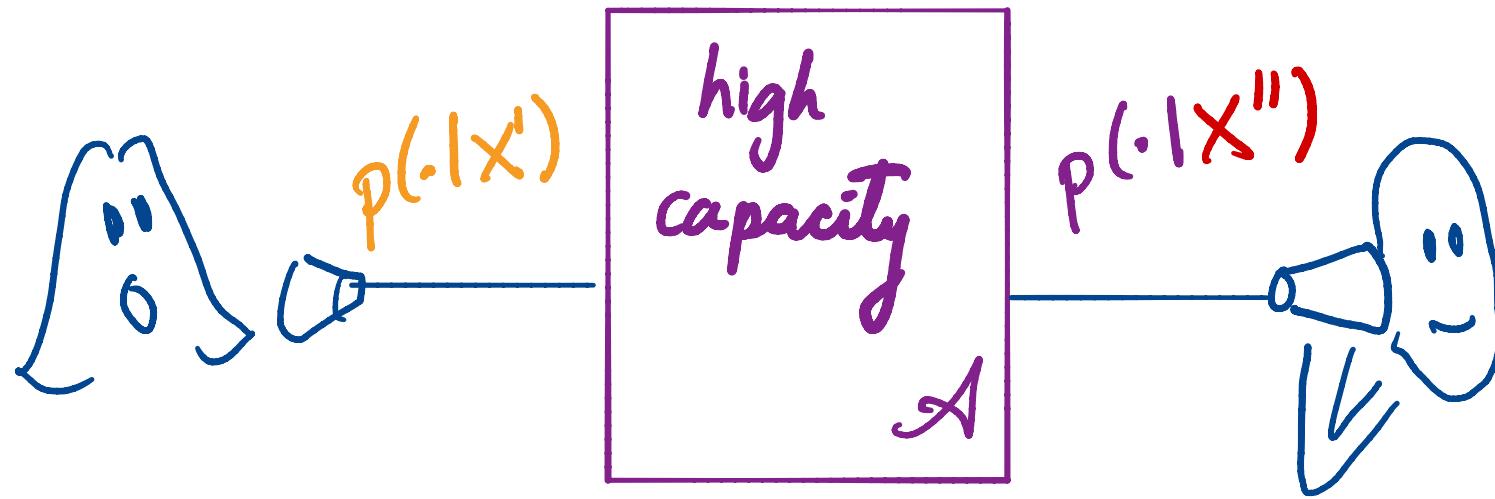
Adjust your algos
to increase exp. log. PA

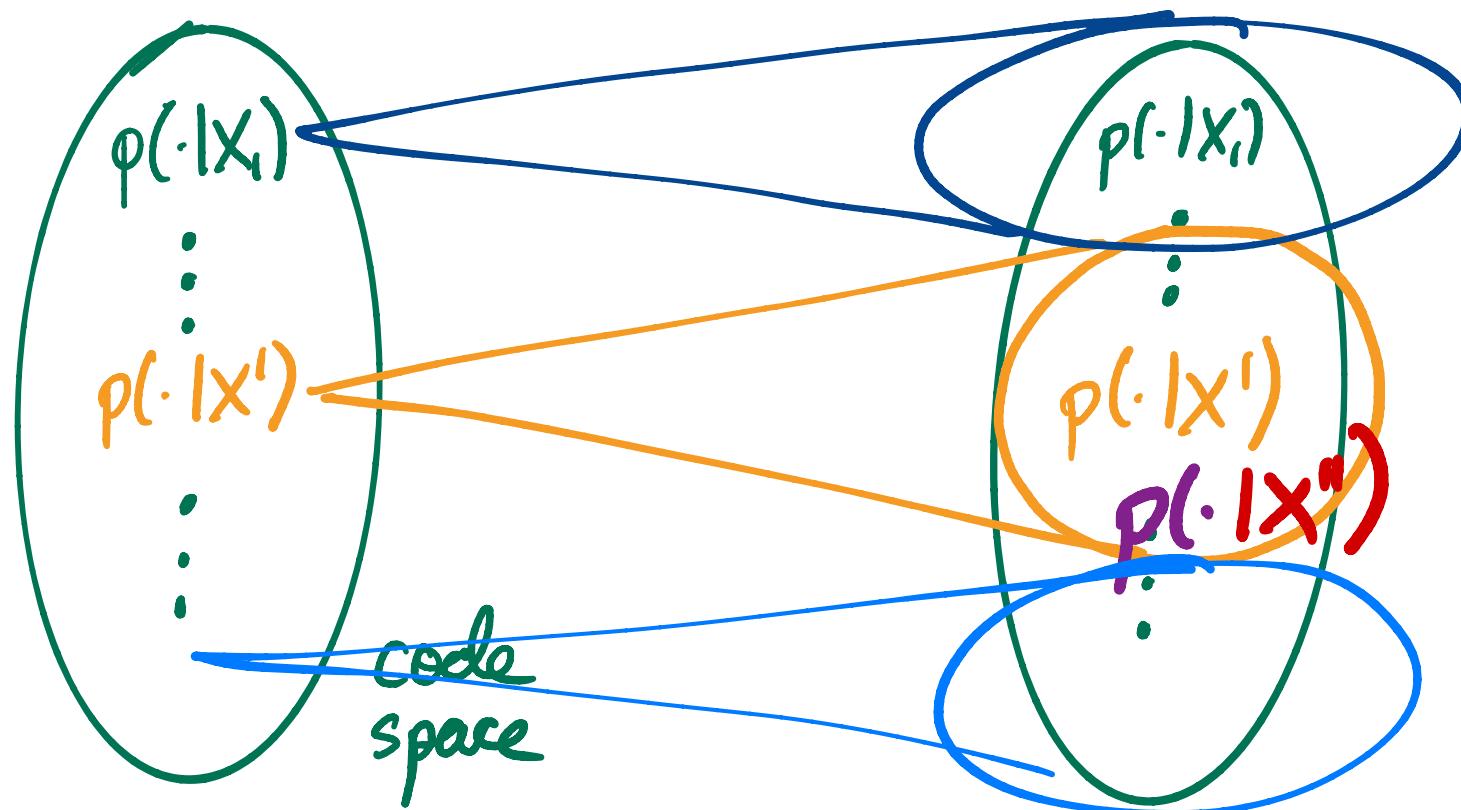
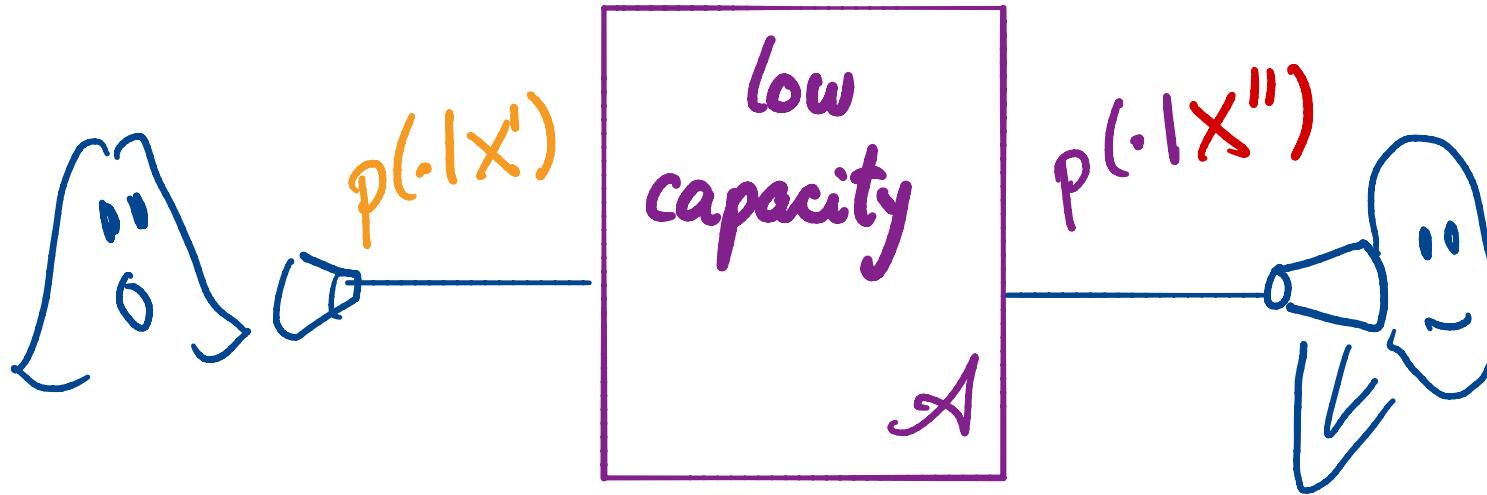
Information transmission

* Codewords have content

* Channel is robust

Adjust your channels
to increase capacity





$$\Pr(E) \leq \text{const} \cdot \exp(-\log |G| (\text{exp. log. PA} - \frac{\log m}{\log |G|} - \varepsilon))$$

if $\text{exp. log. PA} > \frac{\log m}{\log |G|}$, then

$$\Pr(E) \xrightarrow[n \rightarrow \infty]{\text{in prob.}} 0$$

codeword length

channel capacity

code rate

$$\Pr(E) \leq 2^{-n} \left(\text{capacity} - \frac{\log \#\text{msgs}}{n} - 3\varepsilon \right)$$

if $\text{cap} > \frac{\log 2^{\text{nr}}}{n} = r$, then

$$\Pr(E) \xrightarrow[n \rightarrow \infty]{} 0.$$

Posterior agreement

Expected log posterior agreement:

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

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

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

