

# CRF Word Alignment & Noisy Channel Translation

## Machine Translation Lecture 6

**Instructor: Chris Callison-Burch**  
**TAs: Mitchell Stern, Justin Chiu**

**Website: [mt-class.org/penn](http://mt-class.org/penn)**



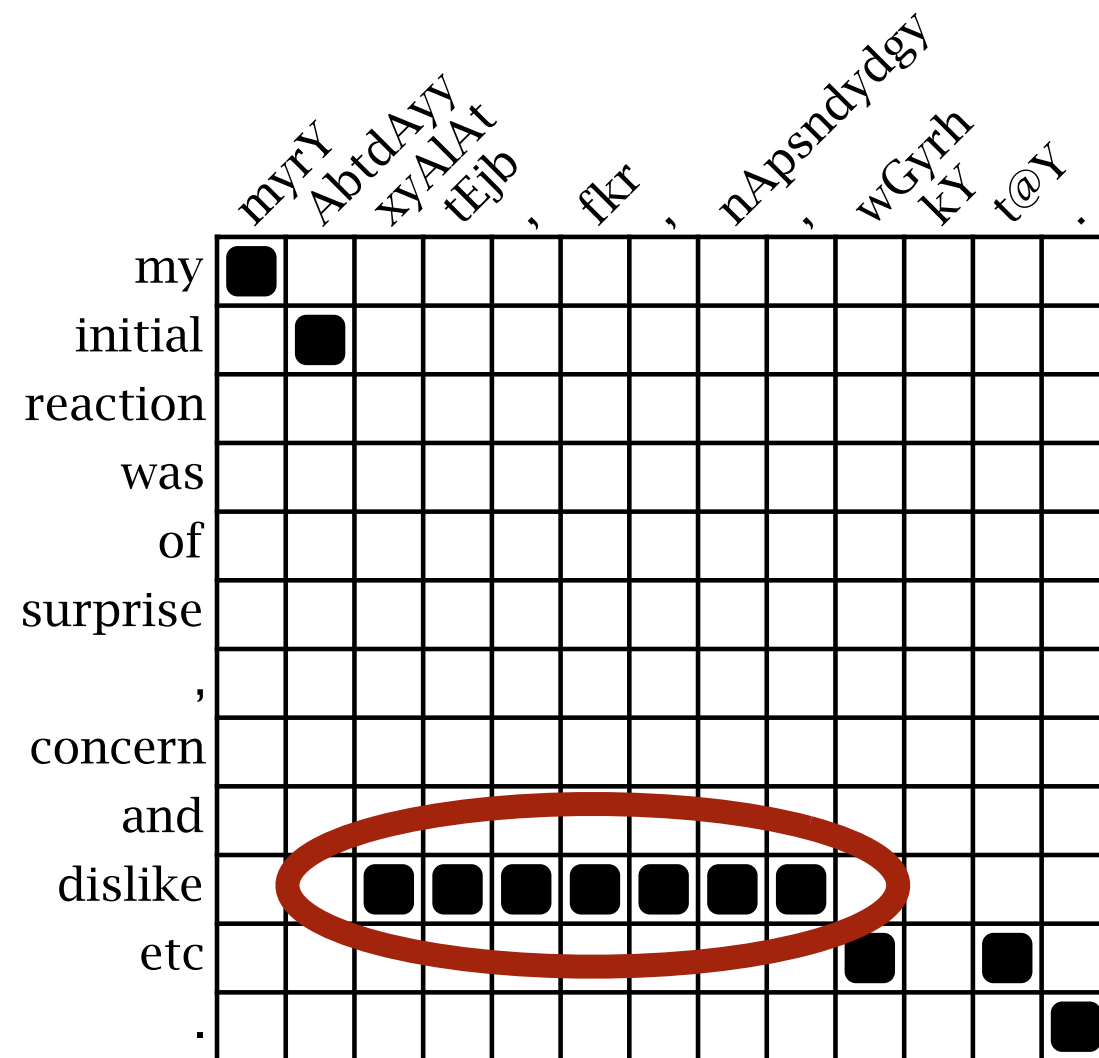
# Last Time ...

$$p(\text{Translation}) = \sum_{\text{Alignment}} p(\text{Alignment}, \text{Translation})$$

$$= \sum_{\text{Alignment}} \underbrace{p(\text{Alignment})}_{\text{Alignment}} \times \underbrace{p(\text{Translation} \mid \text{Alignment})}_{\text{Translation}}$$

$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in [0, n]^m} \underbrace{p(\mathbf{a} \mid \mathbf{f}, m)}_{\text{Alignment}} \times \prod_{i=1}^m \underbrace{p(e_i \mid f_{a_i})}_{\text{Translation}}$$

# MAP alignment



IBM Model 4 alignment

# A few tricks...

$$p(f|e)$$

	michael	geht	davon	aus	,	dass	er	im	haus	bleibt
michael	1									
assumes		1	1	1						
that						1				
he							1			
will										
stay										1
in								1		
the										
house									1	

English to German

# A few tricks...

$p(f|e)$

	michael	geht	davon	aus	.	dass	er	im	haus	bleibt
michael	■									
assumes		■	■	■						
that						■				
he							■			
will										
stay										■
in								■		
the										
house									■	

English to German

	michael	geht	davon	aus	.	dass	er	im	haus	bleibt
michael	■									
assumes		■								
that						■				
he							■			
will										■
stay										
in								■		
the										
house									■	

German to English

$p(e|f)$

# A few tricks...

$p(f|e)$

	michael	geht	davon	aus	.	dass	er	im	haus	bleibt
michael	■									
assumes		■	■	■						
that						■				
he							■			
will										
stay									■	
in							■			
the										
house									■	

English to German

	michael	geht	davon	aus	.	dass	er	im	haus	bleibt
michael	■									
assumes		■								
that						■				
he							■			
will										■
stay										
in							■			
the										
house									■	

German to English

$p(e|f)$

	michael	geht	davon	aus	.	dass	er	im	haus	bleibt
michael	■									
assumes		■	■	■						
that						■				
he							■			
will										■
stay									■	
in							■			
the								■		
house									■	


Intersection / Union

# Another View

With this model:

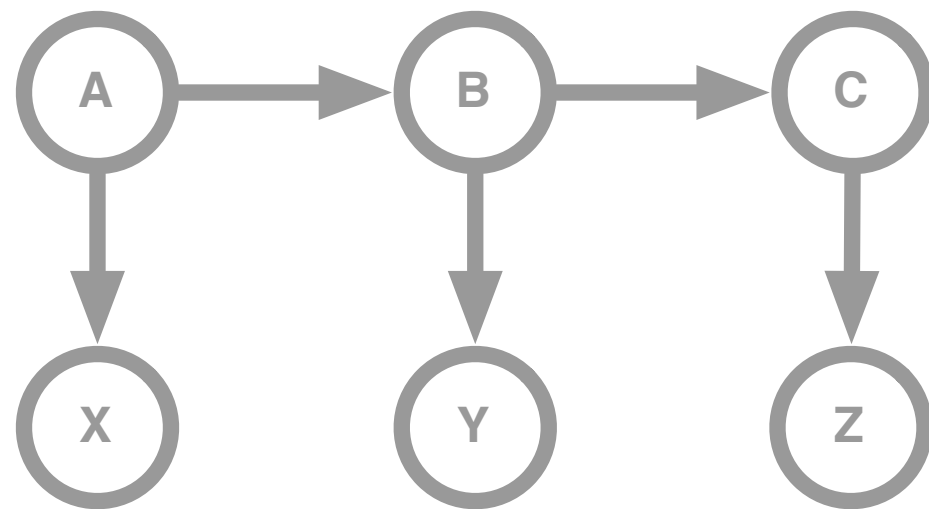
$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in [0, n]^m} p(\mathbf{a} \mid \mathbf{f}, m) \times \prod_{i=1}^m p(e_i \mid f_{a_i})$$

The problem of word alignment is as:

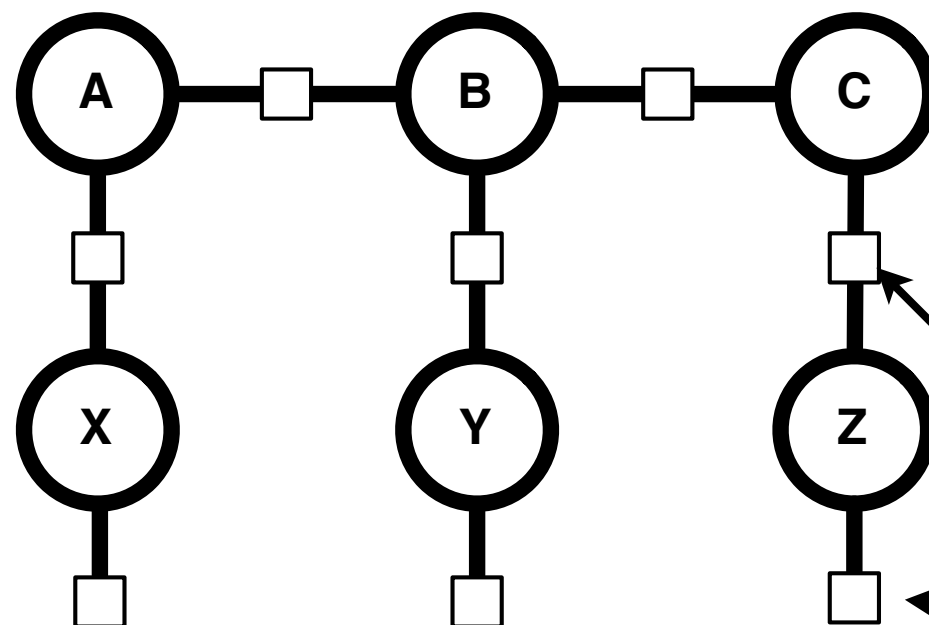
$$\mathbf{a}^* = \arg \max_{\mathbf{a} \in [0, n]^m} p(\mathbf{a} \mid \mathbf{e}, \mathbf{f}, m)$$


Can we model this distribution directly?

# Markov Random Fields (MRFs)



$$p(A, B, C, X, Y, Z) = p(A) \times p(B \mid A) \times p(C \mid B) \times p(X \mid A) p(Y \mid B) p(Z \mid C)$$

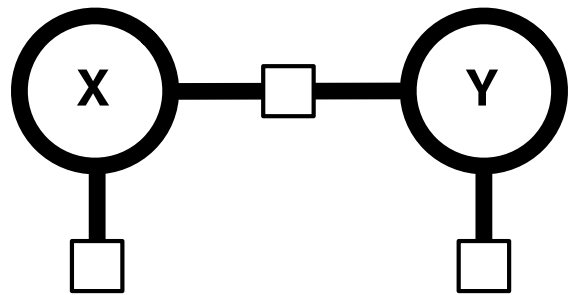


$$p(A, B, C, X, Y, Z) = \frac{1}{Z} \times \Psi_1(A, B) \times \Psi_2(B, C) \times \Psi_3(C, D) \times \Psi_4(X) \times \Psi_5(Y) \times \Psi_6(Z)$$

“Factors”



# Computing $Z$



$$\mathcal{X} = \{a, b, c\}$$

$$X \in \mathcal{X}$$

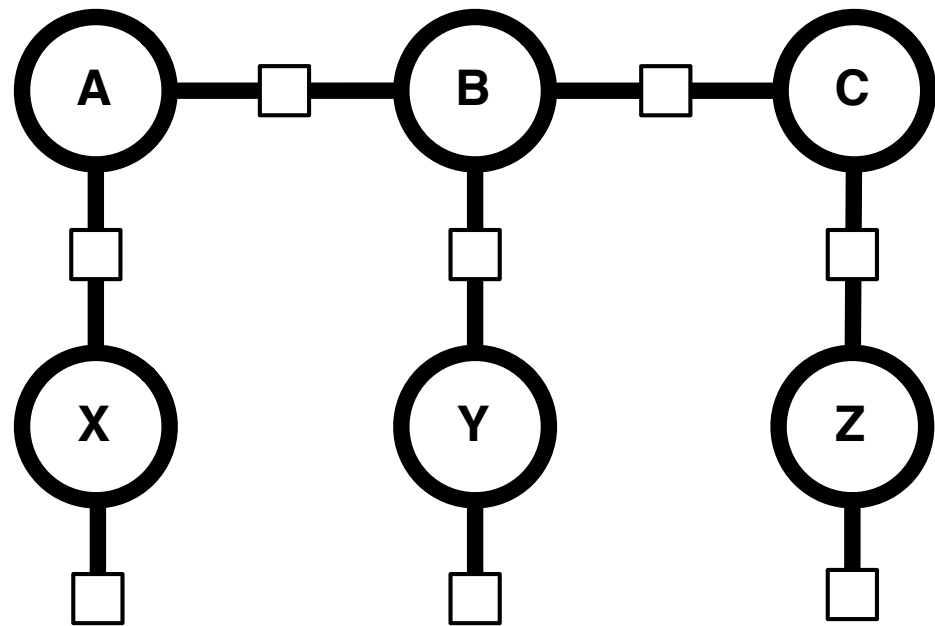
$$Y \in \mathcal{X}$$

$$Z = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{X}} \Psi_1(x, y) \Psi_2(x) \Psi_3(y)$$

When the graph has certain structures (e.g., chains), you can factor to get polynomial time dynamic programming algorithms.

$$Z = \sum_{x \in \mathcal{X}} \Psi_2(x) \sum_{y \in \mathcal{X}} \Psi_1(x, y) \Psi_3(y)$$

# Log-linear models



$$p(A, B, C, X, Y, Z) = \frac{1}{Z} \times \\ \Psi_1(A, B) \times \Psi_2(B, C) \times \Psi_3(C, D) \times \\ \Psi_4(X) \times \Psi_5(Y) \times \Psi_6(Z)$$

$$\Psi_{1,2,3}(x, y) = \exp \sum_k w_k f_k(x, y)$$

Weights (learned)

Feature functions  
(specified)

# Random Fields


- **Benefits**
  - Potential functions can be defined with respect to arbitrary features (functions) of the variables
  - Great way to incorporate knowledge
- **Drawbacks**
  - Likelihood involves computing  $Z$
  - Maximizing likelihood usually requires computing  $Z$  (often over and over again!)

# Conditional Random Fields

- Use MRFs to parameterize a conditional distribution. Very easy: let feature functions look at **anything** they want in the “input”

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z_{\mathbf{w}}(\mathbf{y})} \exp \sum_{F \in \mathcal{G}} \sum_k w_k f_k(F, \mathbf{x})$$

All factors in the graph of  $\mathbf{y}$



# Parameter Learning

- CRFs are trained to maximize conditional likelihood

$$\hat{\mathbf{w}}_{\text{MLE}} = \arg \max_{\mathbf{w}} \prod_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}} p(\mathbf{y}_i \mid \mathbf{x}_i ; \mathbf{w})$$

- Recall we want to directly model

$$p(\mathbf{a} \mid \mathbf{e}, \mathbf{f})$$

- The likelihood of what alignments?

**Gold reference alignments!**

# CRF for Alignment

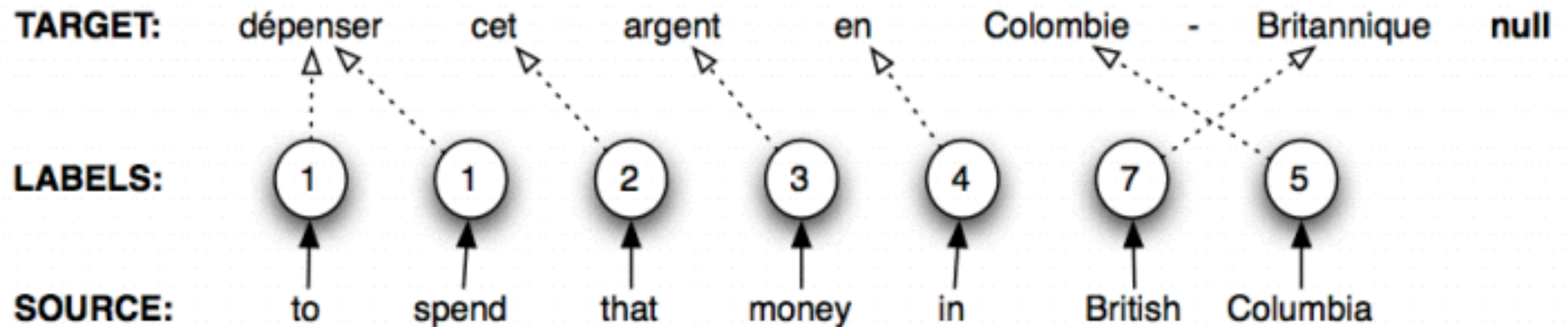
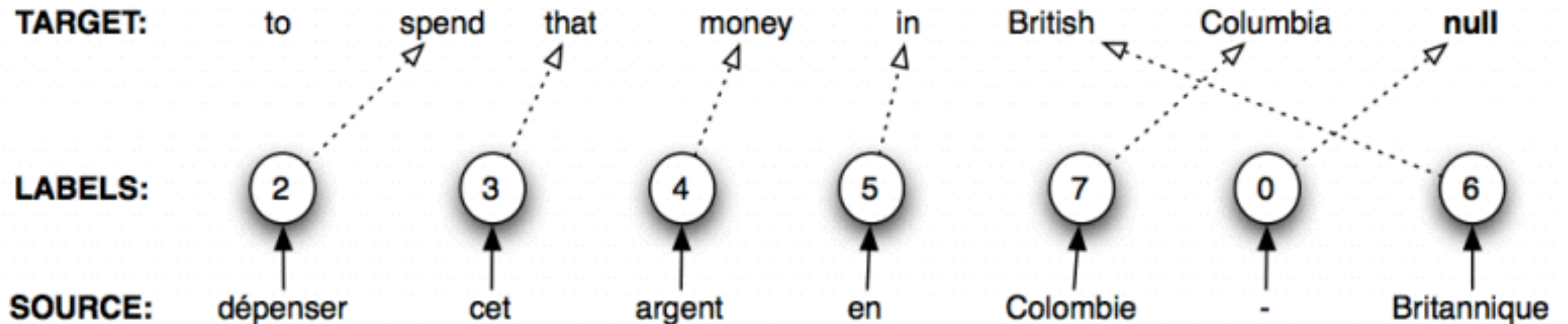
- One of many possibilities, due to Blunsom & Cohn (2006)

$$p(\mathbf{a} \mid \mathbf{e}, \mathbf{f}) = \frac{1}{Z_{\mathbf{w}}(\mathbf{e}, \mathbf{f})} \exp \sum_{i=1}^{|\mathbf{e}|} \sum_k w_k f(a_i, a_{i-1}, i, \mathbf{e}, \mathbf{f})$$

- $\mathbf{a}$  has the same form as in the lexical translation models (still make a one-to-many assumption)
- $w_k$  are the model parameters
- $f_k$  are the feature functions

$$O(n^2 m) \approx O(n^3)$$

# Model

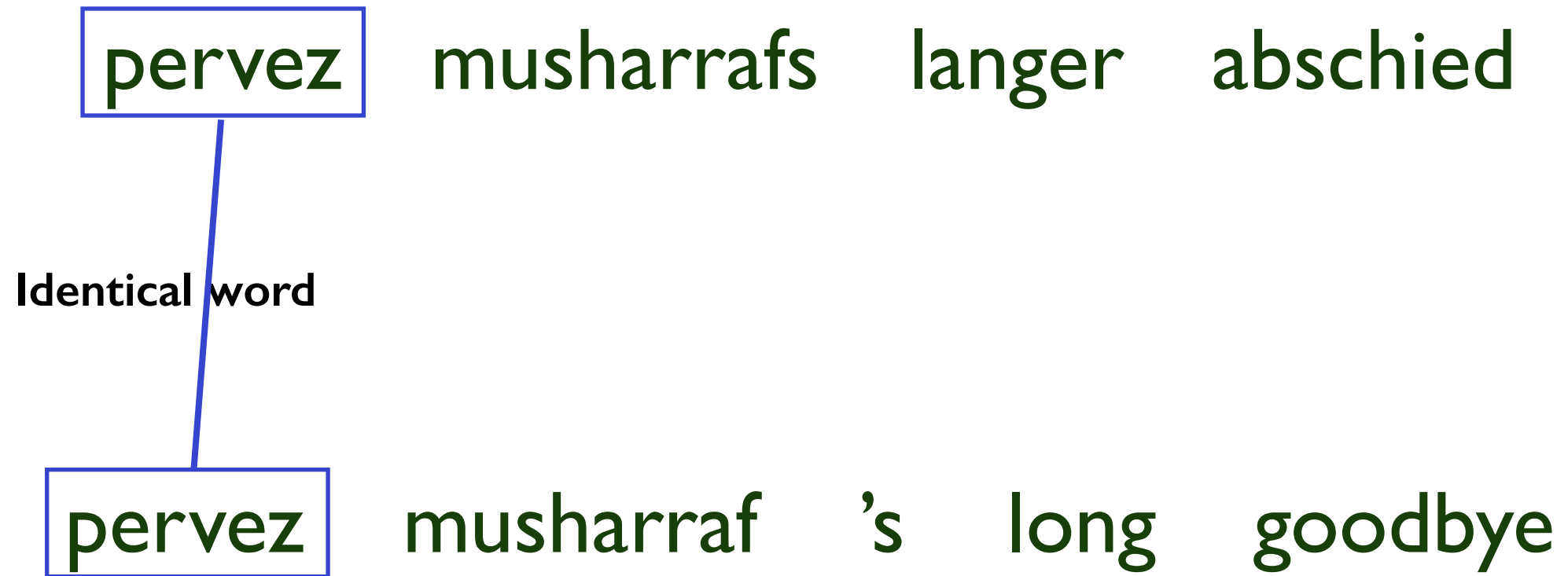


- Labels (one per target word) index source sentence
- Train model (e,f) and (f,e) [inverting the reference alignments]

# Alignment Experiments

- French-English Canadian Hansards corpus
- 484 manually word-aligned sentence pairs  
(100 training, 37 development, 347 testing)
- 1.1 million sentence-aligned pairs
- Baseline for comparison: Giza++  
implementation of IBM Model 4
- (Also experimented on Romanian-English)





Identical word

pervez **musharrafs** langer abschied

Matching prefix

pervez **musharra**f 's long goodbye

**Identical word**  
**Matching prefix**

pervez **musharrafs** langer abschied

Matching suffix

pervez musharra**f** **'s** long goodbye

**Identical word**

**Matching prefix**

**Matching suffix**

pervez musharrafs **langer** abschied

Orthographic similarity

pervez musharrafa 's **long** goodbye

**Identical word**

**Matching prefix**

**Matching suffix**

**Orthographic similarity**

pervez musharrafs langer **abschied**

In dictionary

pervez musharraf 's long **goodbye**

**Identical word**

**Matching prefix**

**Matching suffix**

**Orthographic similarity**

**In dictionary**

...

# Lexical Features

- Word $\leftrightarrow$ word indicator features
- Various word $\leftrightarrow$ word co-occurrence scores
  - IBM Model 1 probabilities ( $t \rightarrow s$  ,  $s \rightarrow t$ )
  - Geometric mean of Model 1 probabilities
  - Dice's coefficient [binned]
  - Products of the above

# Lexical Features

- Word class ↔ word class indicator
  - **NN** translates as **NN** (NN\_NN=1)
  - **NN** does not translate as **MD** (NN\_MD=1)
- Identical word feature
  - **2010 = 2010** (IdentWord=1 IdentNum=1)
- Identical prefix feature
  - **Obama**<sub>a</sub> ~ **Obamu**<sub>u</sub> (IdentPrefix=1)
- Orthographic similarity measure [binned]
  - **Al-Qa**<sub>a</sub>**eda** ~ **Al-Ka**<sub>a</sub>**ida** (OrthoSim050\_080=1)

# Other Features

- Compute features from large amounts of unlabeled text
- Does the Model 4 alignment contain this alignment point?
- What is the Model 1 posterior probability of this alignment point?



# Results

Alignment Results:

	Precision	Recall	F-score
French $\rightarrow$ English	0.97	0.86	0.91
French $\leftarrow$ English	<b>0.98</b>	0.83	0.91
French $\leftrightarrow$ English	0.96	0.90	0.93
French $\rightarrow$ English (+ibm model4)	<b>0.98</b>	0.88	0.93
French $\leftarrow$ English (+ibm model4)	<b>0.98</b>	0.87	0.93
French $\leftrightarrow$ English (+ibm model4)	<b>0.98</b>	0.91	<b>0.95</b>
GIZA++ (French $\leftrightarrow$ English)	0.87	<b>0.95</b>	0.91

# Summary

- CRFs outperform unsupervised / latent variable alignment models, even when only a small number of word-aligned sentences are available
- Diverse range of features can be incorporated and are beneficial to word-alignment quality
- Features from unsupervised models can also be incorporated

Unfortunately, you need gold alignments!

# Putting the pieces together

- We have seen how to model the following:

$$p(\mathbf{e})$$

$$p(\mathbf{e} \mid \mathbf{f}, m)$$

$$p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m)$$

$$p(\mathbf{a} \mid \mathbf{e}, \mathbf{f})$$

# Putting the pieces together

- We have seen how to model the following:

$$p(\mathbf{e})$$

$$p(\mathbf{e} \mid \mathbf{f}, m)$$

$$p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m)$$

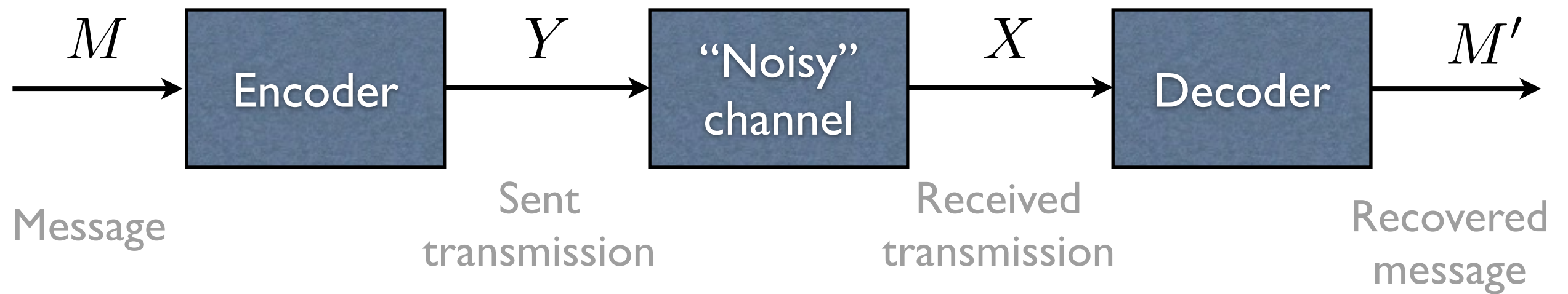
$$p(\mathbf{a} \mid \mathbf{e}, \mathbf{f})$$

- Goal: a better model of  $p(\mathbf{e} \mid \mathbf{f}, m)$  that knows about  $p(\mathbf{e})$

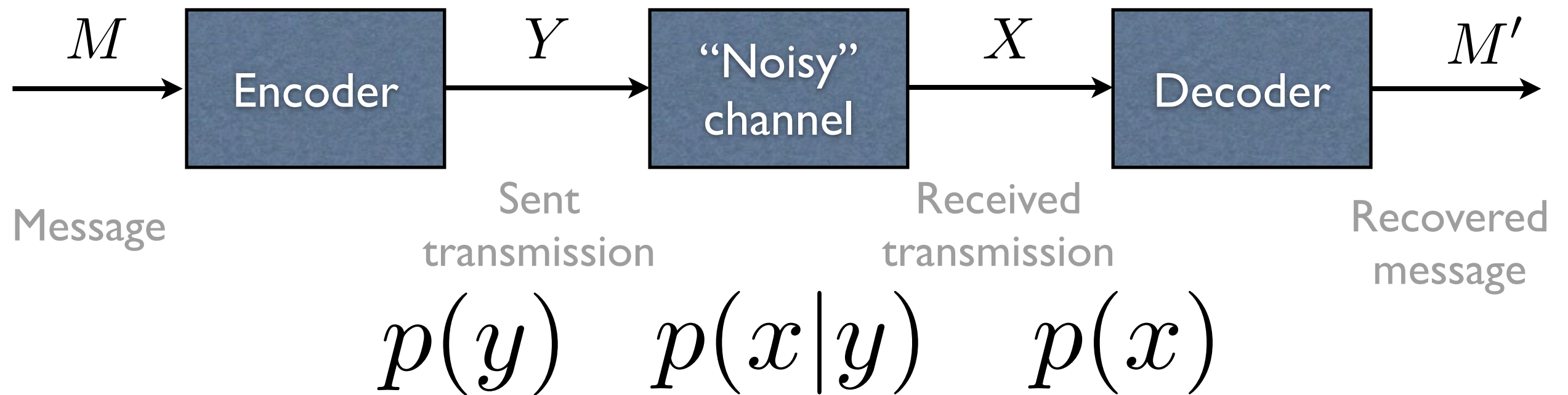
One naturally wonders if the problem of translation could conceivably be treated as a problem in cryptography. When I look at an article in Russian, I say: *‘This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.’*



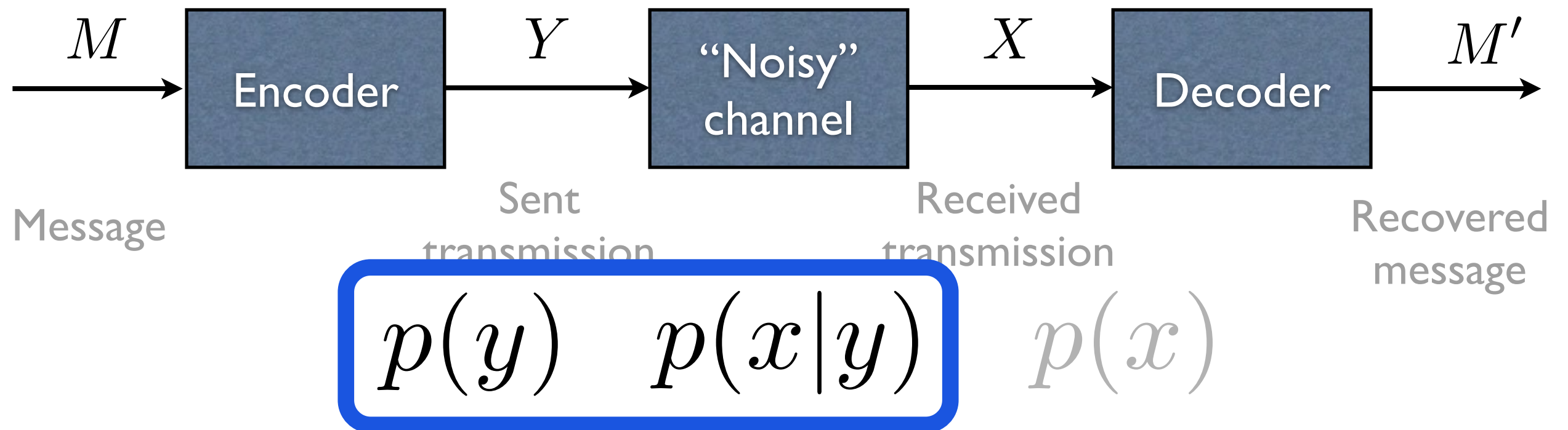
Warren Weaver to Norbert Wiener, March, 1947



Claude Shannon. "A Mathematical Theory of Communication" 1948.

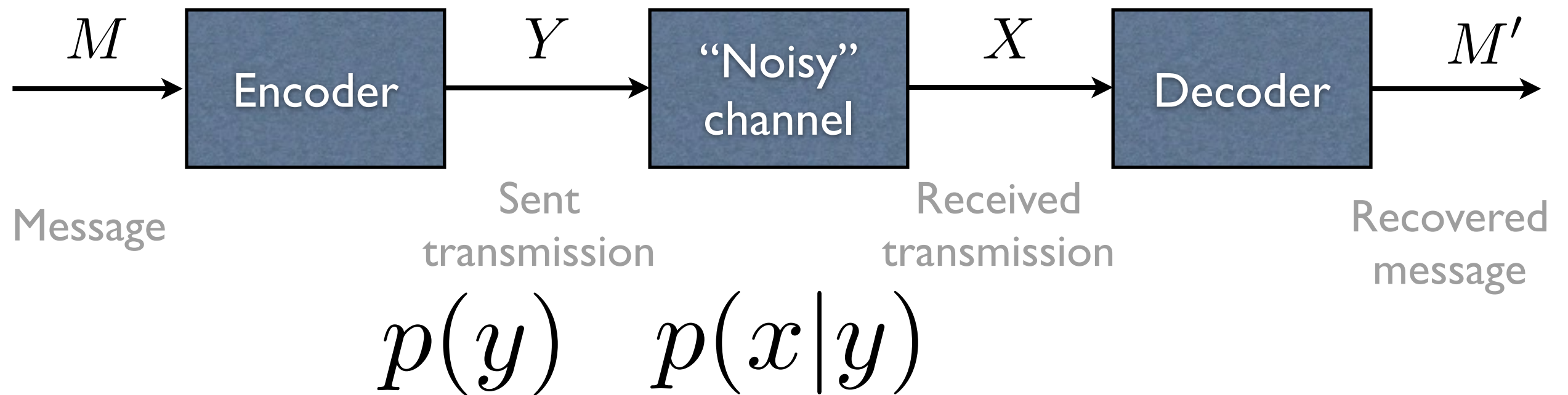


Claude Shannon. "A Mathematical Theory of Communication" 1948.



Claude Shannon. "A Mathematical Theory of Communication" 1948.



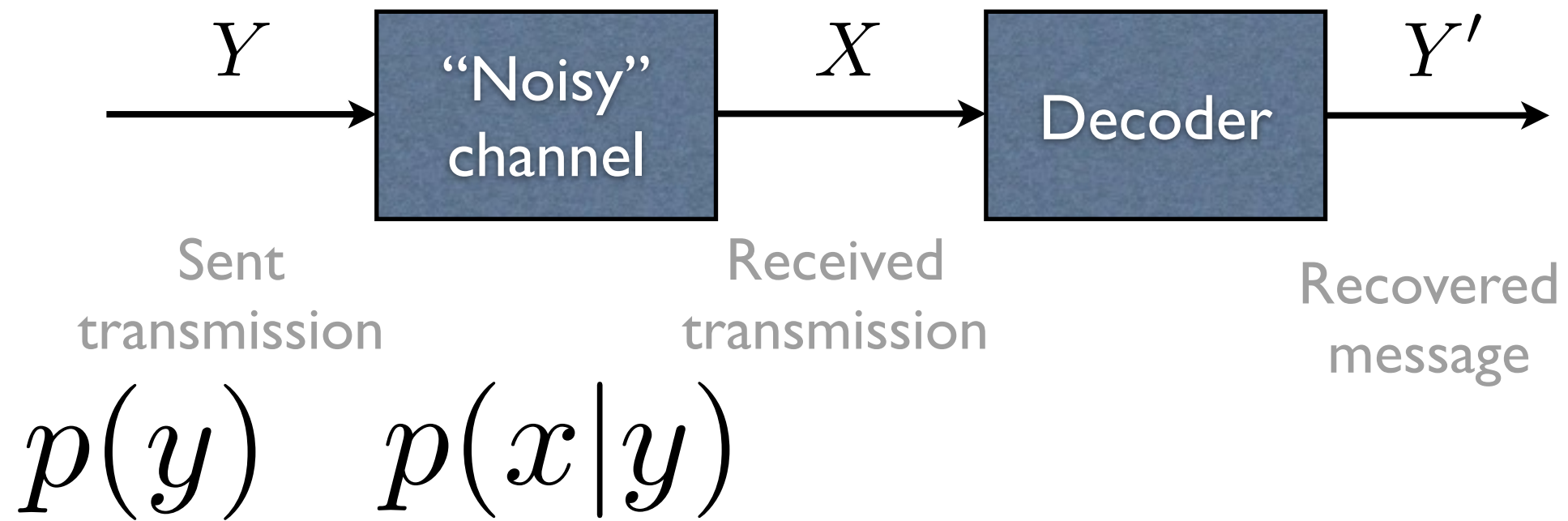


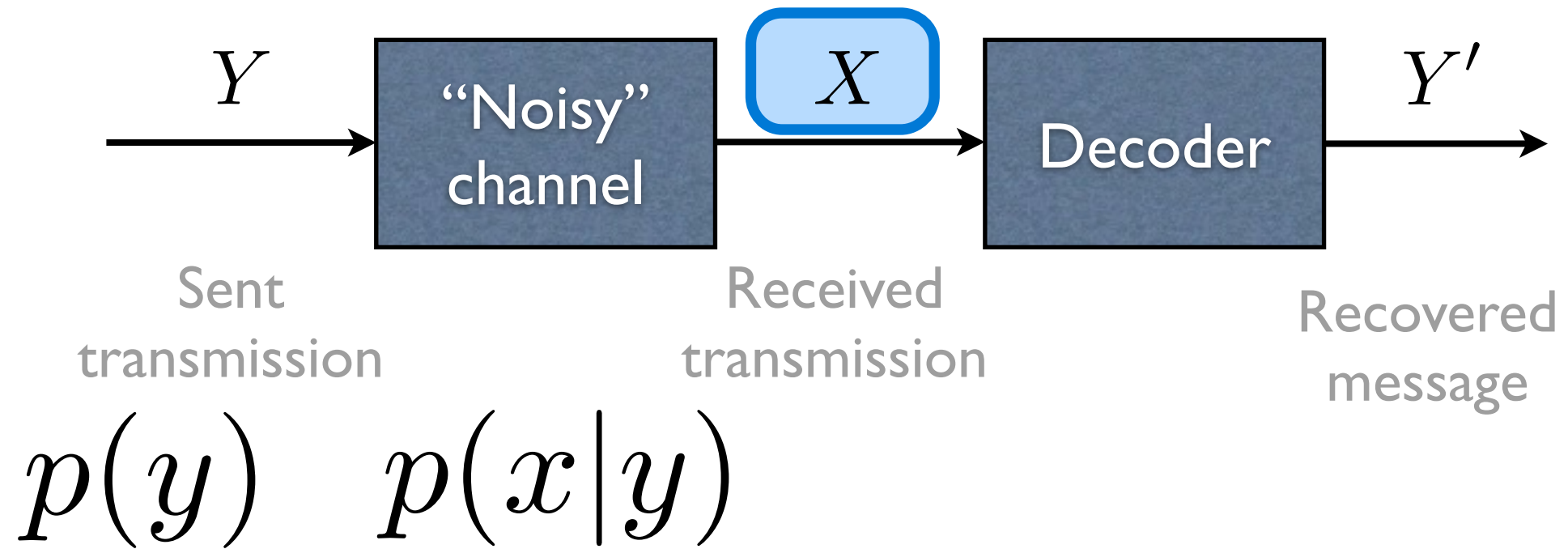
## Shannon's theory tells us:

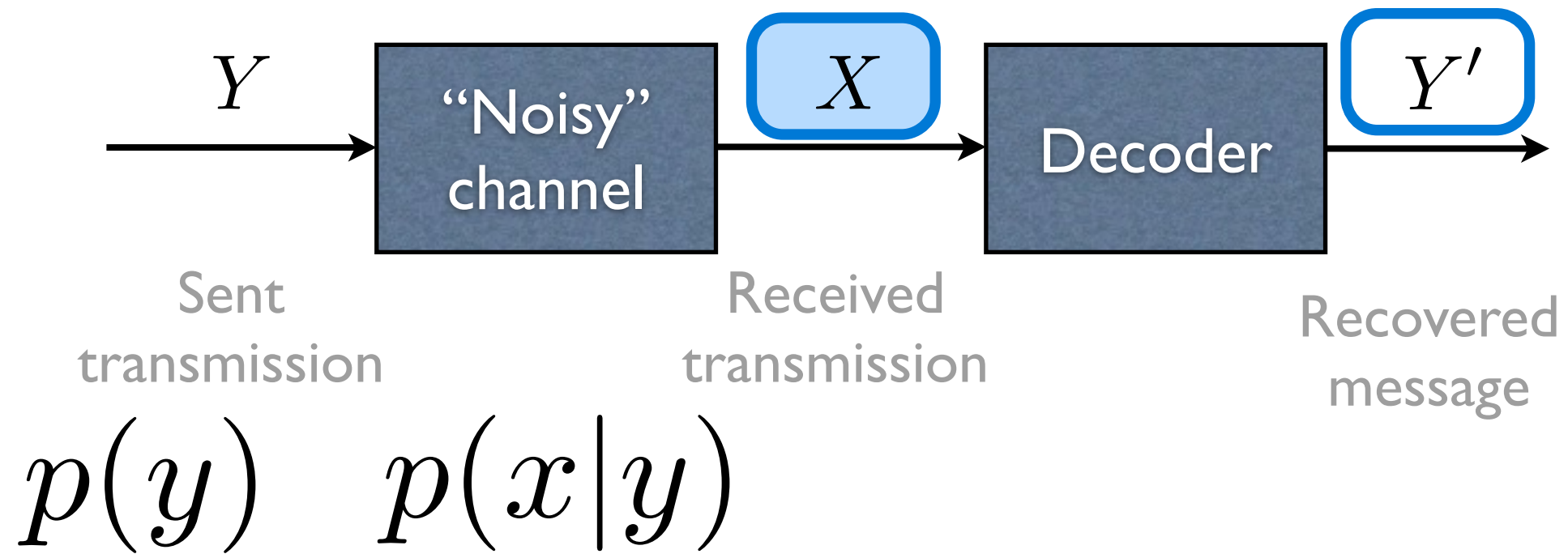
- 1) how much data you can send
- 2) the limits of compression
- 3) why your download is so slow
- 4) how to translate

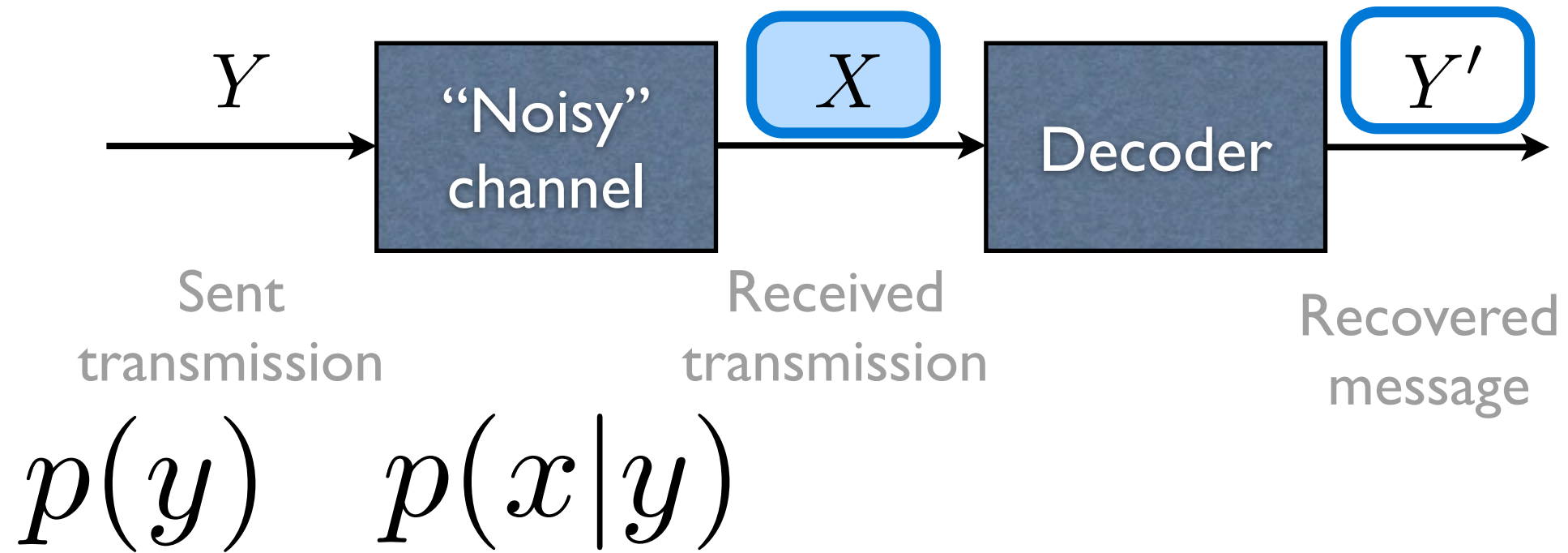


Claude Shannon. "A Mathematical Theory of Communication" 1948.

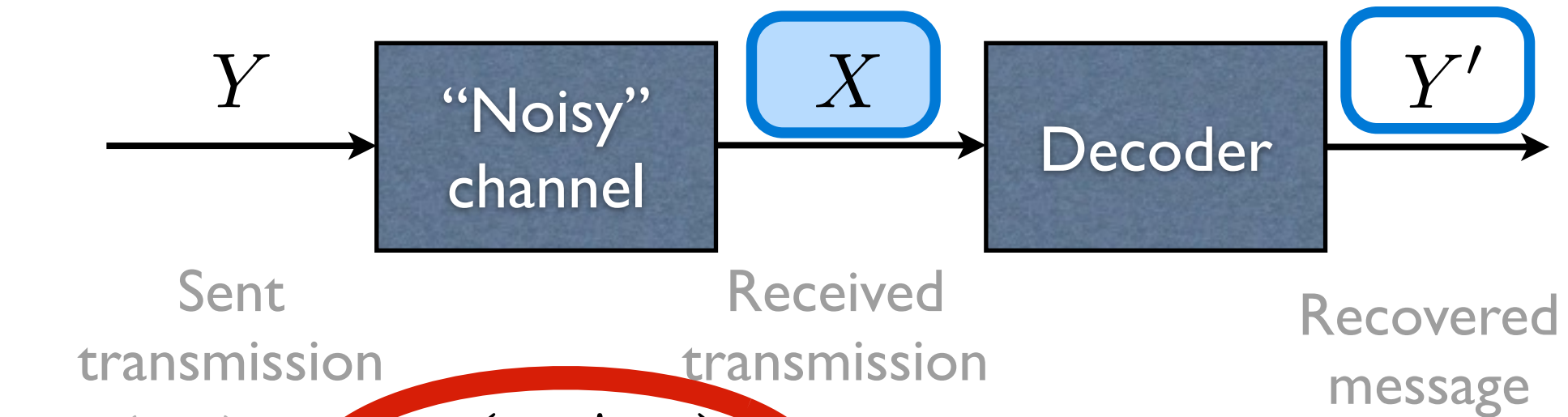








$$\boxed{y'} = \arg \max_y p(y|x)$$



$p(y)$

$$p(x|y)$$

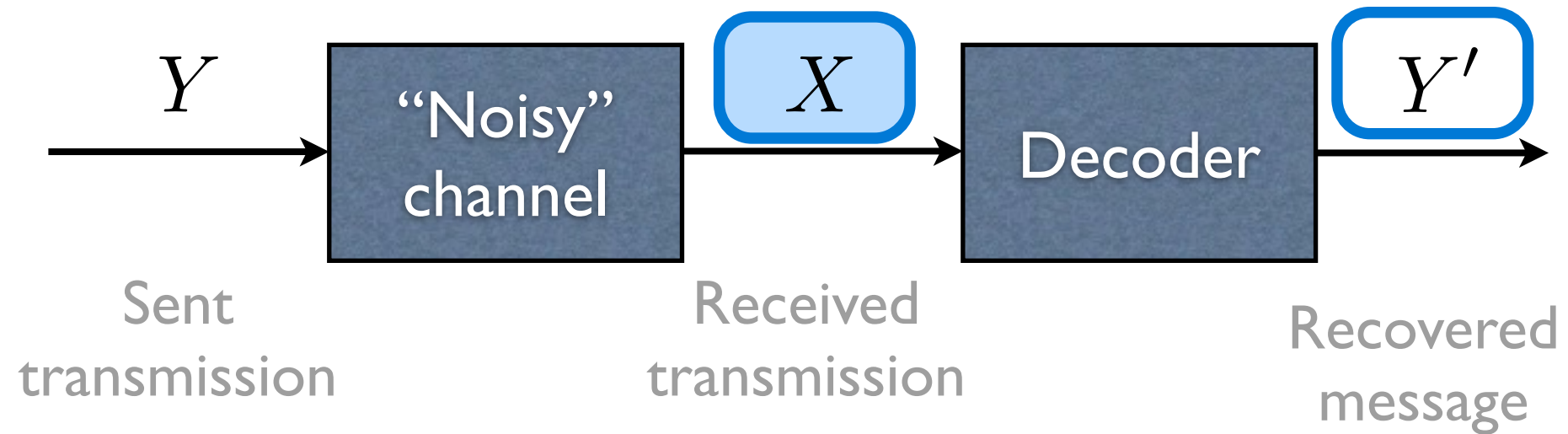
$\neq$

$$y'$$

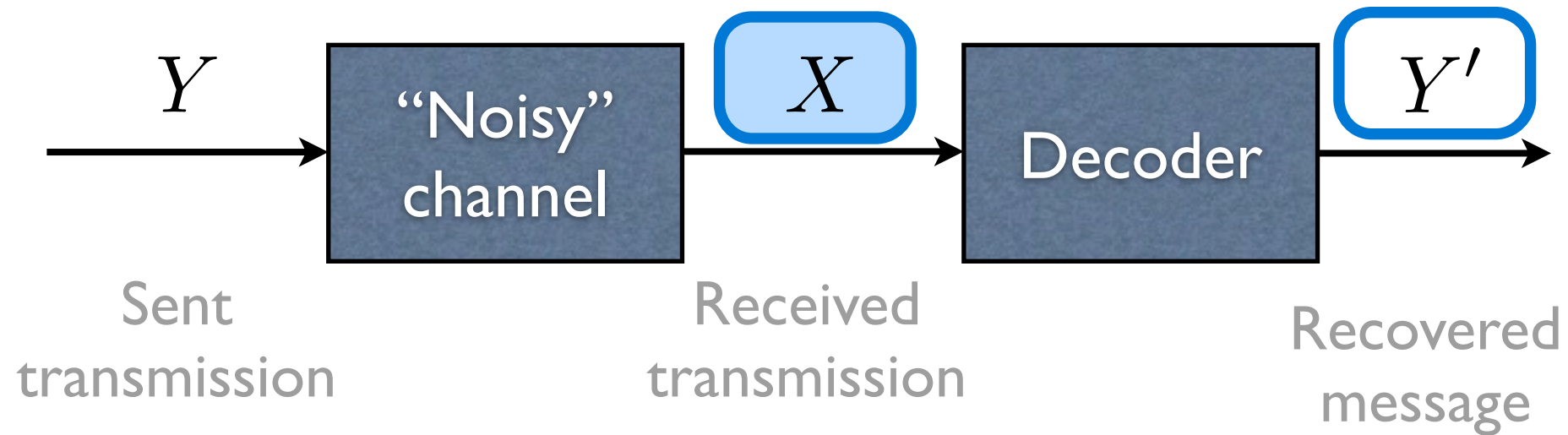
$$= \arg \max_y p(y|x)$$



I can help.



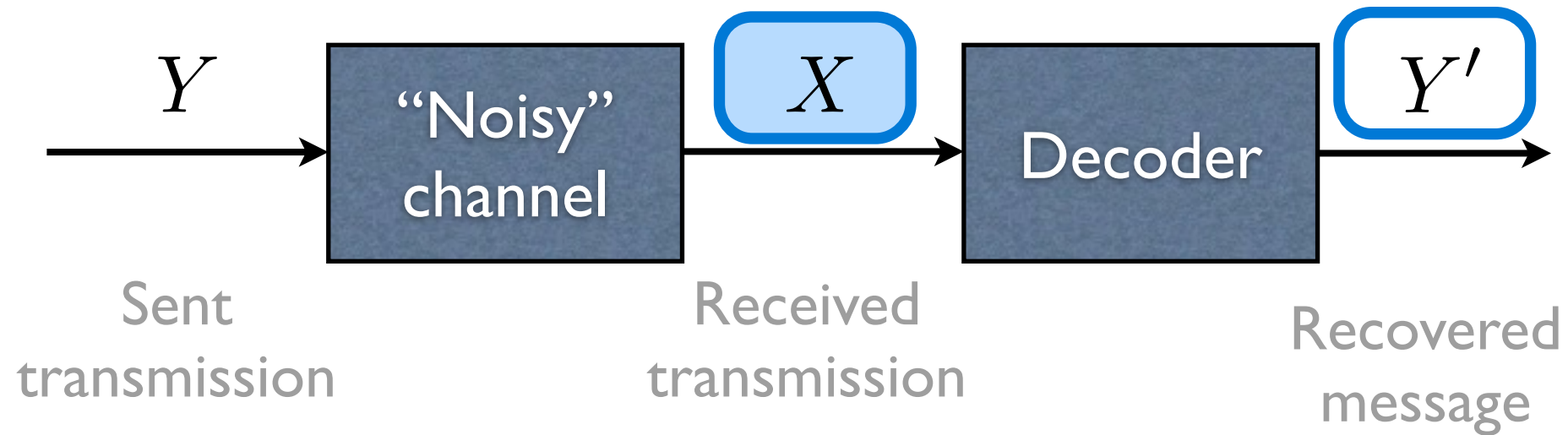
$$\boxed{y'} = \arg \max_y p(y|x)$$
$$= \arg \max_y \frac{p(x|y)p(y)}{p(x)}$$



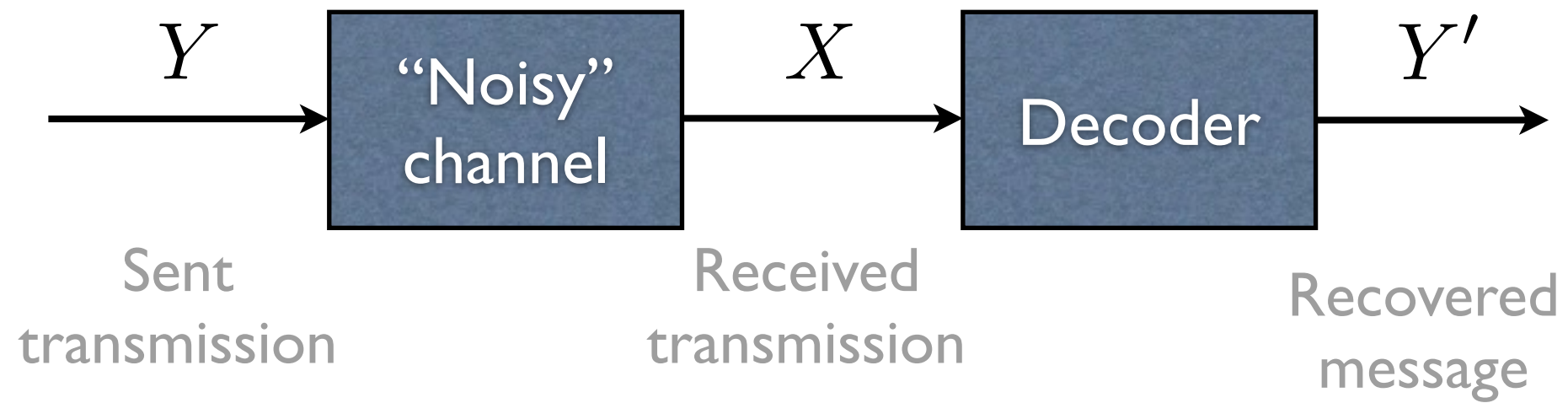
$$\boxed{y'} = \arg \max_y p(y|x)$$
$$= \arg \max_y \frac{p(x|y)p(y)}{p(x)}$$

Denominator doesn't depend on  $y$ .

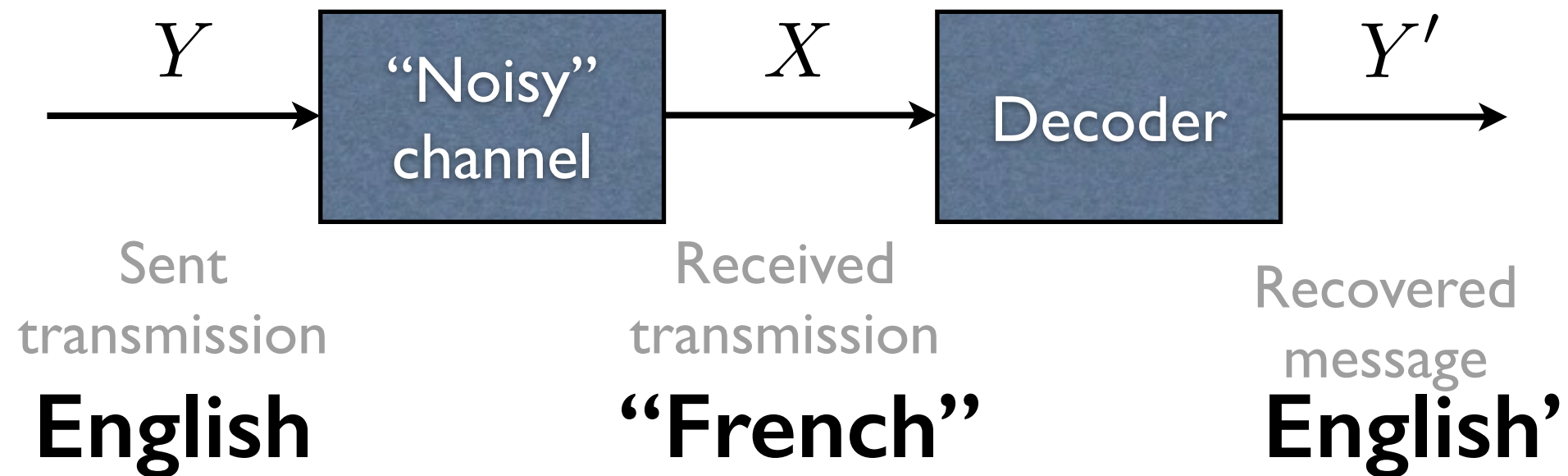




$$\begin{aligned} \boxed{y'} &= \arg \max_y p(y|x) \\ &= \arg \max_y \frac{p(x|y)p(y)}{p(x)} \\ &= \arg \max_y p(x|y)p(y) \end{aligned}$$



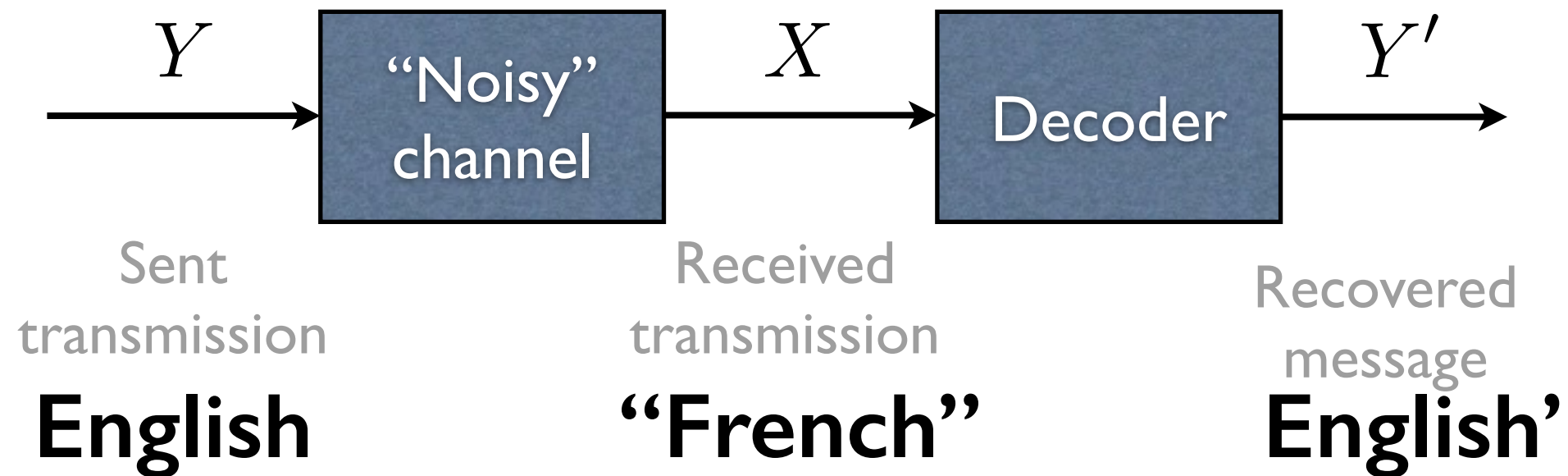
$$y' = \arg \max_y p(x|y)p(y)$$



---

$$y' = \arg \max_y p(x|y)p(y)$$

$$\mathbf{e}' = \arg \max_{\mathbf{e}} p(\mathbf{f}|\mathbf{e})p(\mathbf{e})$$



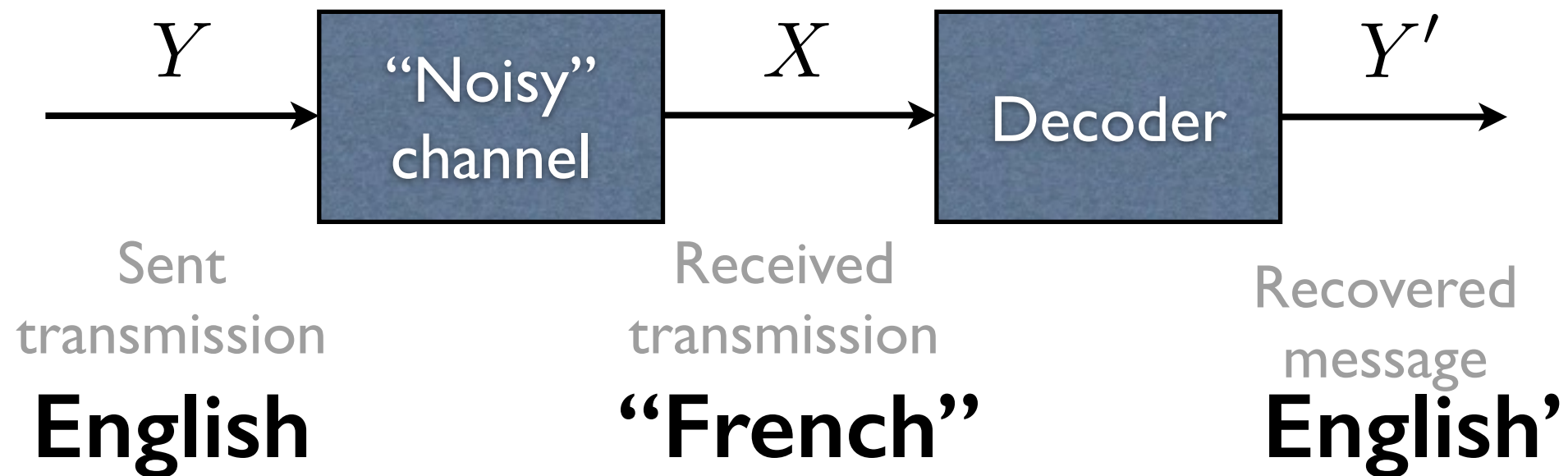
---

$$y' = \arg \max_y p(x|y)p(y)$$

$$e' = \arg \max_e p(\mathbf{f}|\mathbf{e})p(\mathbf{e})$$



translation model



---

$$y' = \arg \max_y p(x|y)p(y)$$

$$e' = \arg \max_e p(\mathbf{f}|\mathbf{e})p(\mathbf{e})$$



translation model



language model

**Other noisy channel applications: OCR, speech recognition, spelling correction...**

# Division of labor

- Translation model
  - probability of translation *back* into the source
  - ensures adequacy of translation
- Language model
  - is a translation hypothesis “good” English?
  - ensures fluency of translation



$$\begin{aligned}
 \mathbf{e}^* &= \arg \max_{\mathbf{e}} p(\mathbf{e} \mid \mathbf{f}) \\
 &= \arg \max_{\mathbf{e}} p(\mathbf{f} \mid \mathbf{e}) \times p(\mathbf{e})
 \end{aligned}$$

# Announcements

- HW1 leaderboard submissions are due tonight at 11:59pm
- HW1 writeup and code are due 24 hours later
- Next week: Phrase-based Machine Translation