



# CS3236 INTRODUCTION TO INFORMATION THEORY

## Lecture 4: Lossless compression

Course given by Pierre Senellart

Material by Stephanie Wehner, with additions by P. Senellart

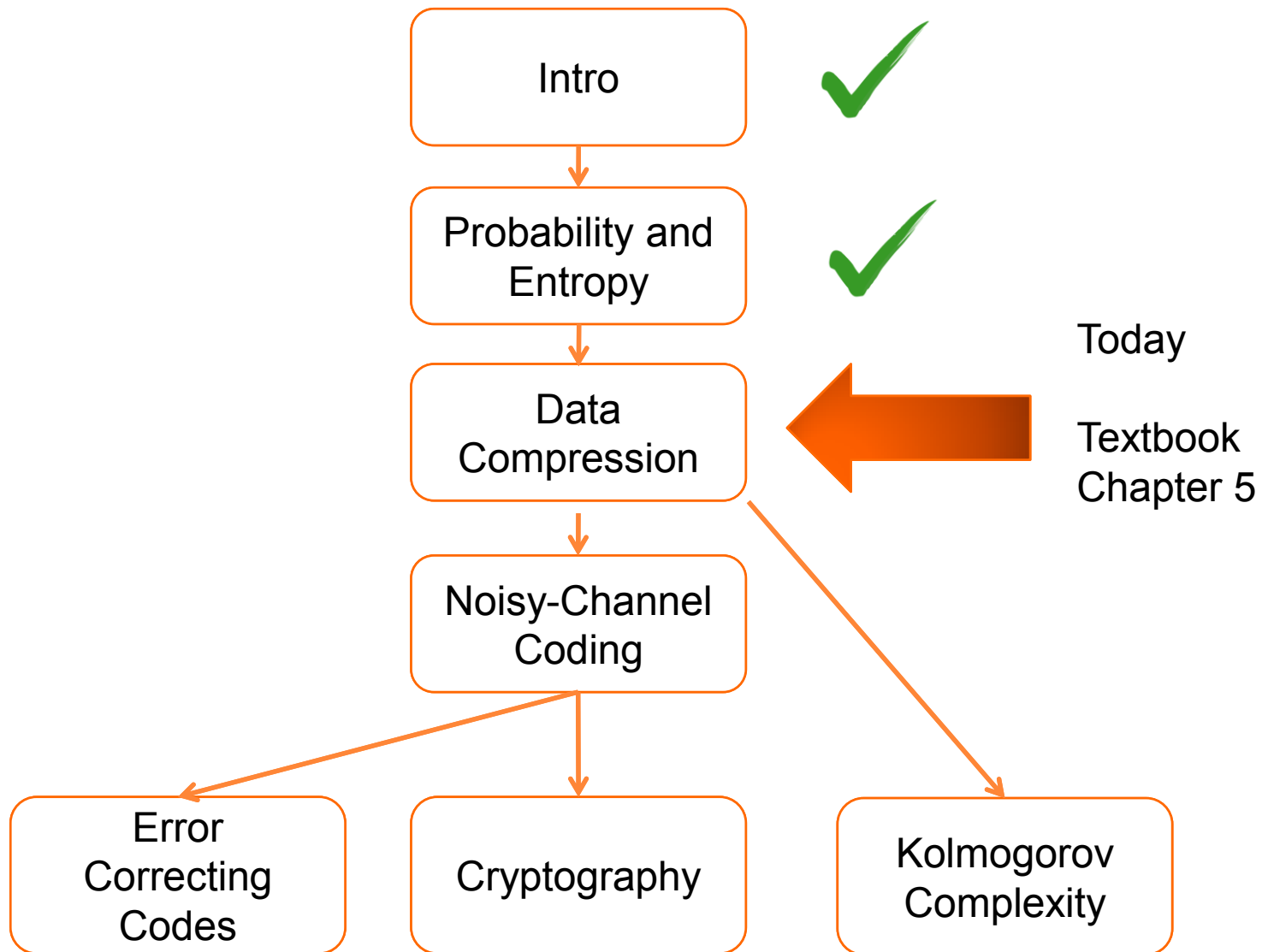
LET'S RECAP...



# CLASS SURVIVAL GUIDE

- Lecture on Monday 2pm-3:30pm
- Tutorial on Wed 4-5pm and Thu 2-3pm, COM1-02-18
  - Choose one session, no need to go to both
  - Discussion, QA, help with homework exercises, projects
- Recess week September 20-28
- No physical lecture/tutorial on Sep 8-12, e-Learning material
- No lecture on October 6 (Hari Raya Haji)
- No tutorial on week 10, October 22-23 (Deepavali)
- Available to meet for answering individual questions. Use <http://www.doodle.com/psenellart> to schedule a meeting. Office Icube #03-09.
- Grading
  - 50% Exam (Friday November 28)
  - 50% Continuous Assessment
    - 10% Mid-Term exam (October 13)
    - 20% Homework (assigned each week, due next Monday)
    - 20% Small project (handed out September 1, due November 14)
- If you have any doubt about the schedule, check IVLE

# WHERE DO WE GO FROM HERE?



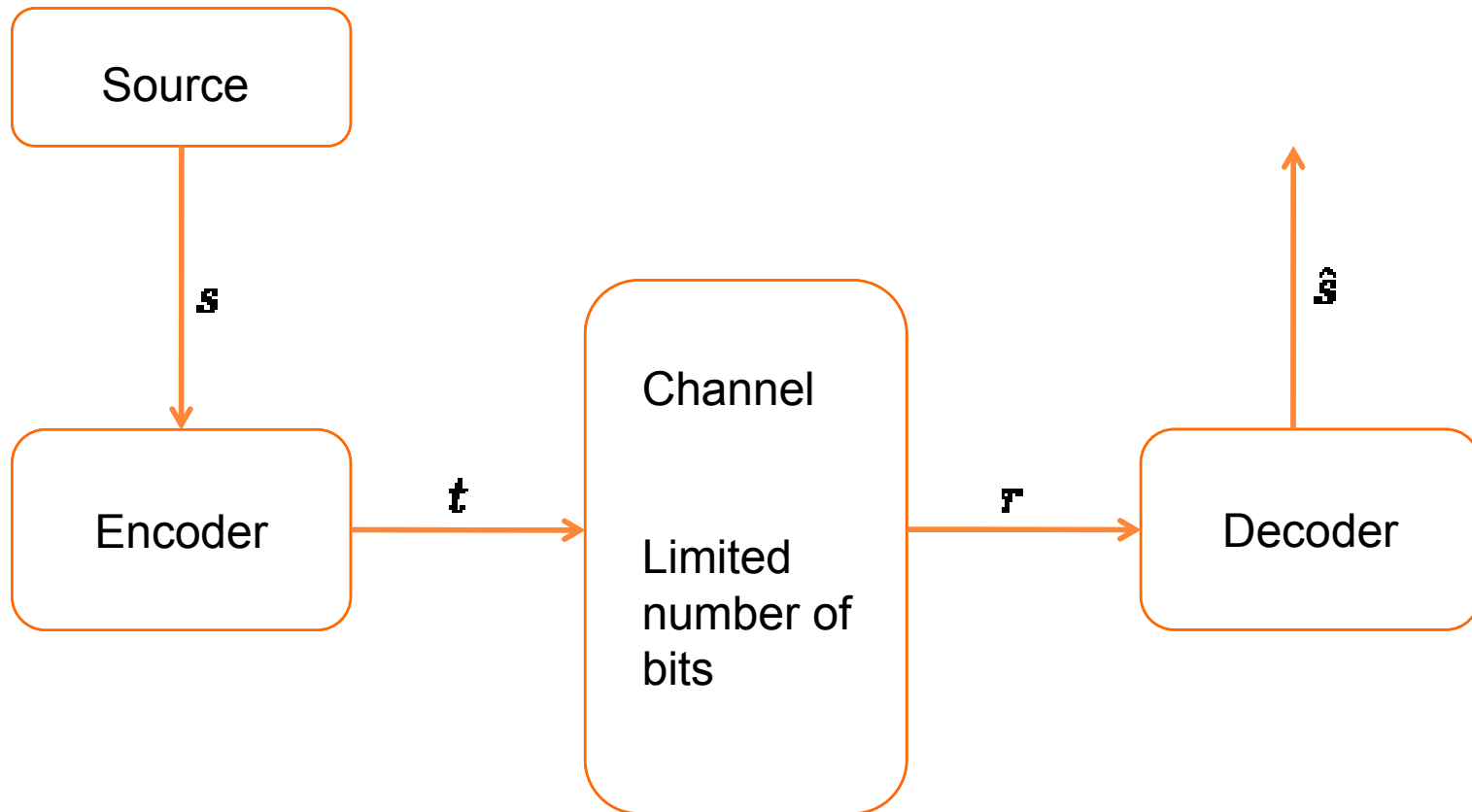
# WHAT WE DID LAST TIME

- Intro to data compression
- Lossy compression framework
- Statement of Shannon's source coding theorem for lossy compression
- The entropy of the source tells us exactly how well we can compress it

# WHAT WE'LL DO TODAY

- Lossless data compression
  - Basic concepts
  - Symbol codes
  - Kraft's inequality
  - Lossless compression and entropy
  - Shannon's source coding theorem for symbol codes
- An optimal symbol code: Huffman coding
- Limits of symbol codes

# DATA COMPRESSION



Goal: Construct an encoder and decoder such that we can recover the information  $\hat{s} = s$

# LOSSY AND LOSSLESS COMPRESSION

## ○ Lossy

- maps some source symbols to the same encoding
- failure probability  $\delta$

## ○ Lossless compression

- maps all possible source symbols to distinct encodings
- works without failure probability

# LOSSLESS COMPRESSION

- If we don't allow a failure probability, how can it be that we can compress at all?
- Idea
  - Use short encodings for symbols that are very likely
  - Use long encodings for symbols that are unlikely
- If we compress many times we will on average compress on average as long encodings are not very frequent

# QUESTIONS IN LOSSLESS COMPRESSION

- Which codewords do we make shorter, and which ones longer?
- How can we make sure that our code is easy to decode in the end?
- What is the best compression that we can achieve?

# SOME TERMINOLOGY: SYMBOL CODE

- A symbol code for an ensemble is a mapping

$$\mathcal{A}_X \rightarrow \{0, 1\}^+$$

- Codeword  $c(x)$
- Length  $l(x)$

All the bit strings of length at least 1

- Example

$$\mathcal{A}_X = \{a, b, c, d\}$$
$$P_X = \left\{ \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8} \right\}$$

x	c(x)	l(x)
a	1000	4
b	0100	4
c	0010	4
d	0001	4

# EXTENDED CODE

- The extended code  $C^+$  is what we get when we use our code to encode an entire string

$$\mathcal{A}_X = \{a, b, c, d\}$$
$$P_X = \left\{ \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8} \right\}$$

$x$	$c(x)$	$l(x)$
a	1000	4
b	0100	4
c	0010	4
d	0001	4

- Example

$$c^+(accd) = 1000001000100001$$

# UNIQUELY DECODABLE

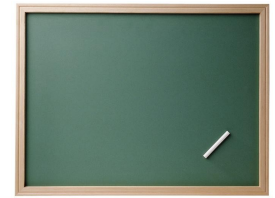
- A code is uniquely decodable if also for the extended code no two strings have the same encoding:

$$\forall x, y \in \mathcal{A}_X^+ \quad x \neq y \implies c^+(x) \neq c^+(y)$$

# REQUIREMENTS OF USEFUL SYMBOL CODES

- The extended code must be uniquely decodable
  - Any string must have a unique decoding
- Must be easy to decode
- Achieve the highest possible compression

# EASY TO DECODE?



- Ideally we want to be able to decode while receiving the signal without storing the signal and looking ahead
- How about this code?

$x$	$c(x)$
a	0
b	011

# PREFIX CODE

- A code is a prefix code if no codeword is a prefix of another

- Not a prefix code

<b>x</b>	<b>c(x)</b>
<b>a</b>	<b>0</b>
<b>b</b>	<b>011</b>

- Prefix code example

<b>x</b>	<b>c(x)</b>
<b>a</b>	<b>1</b>
<b>b</b>	<b>011</b>

- Does this help?

# PREFIX CODES AS TREES

- Let's consider some examples

0 → 0

1 → 10

2 → 111

3 → 110

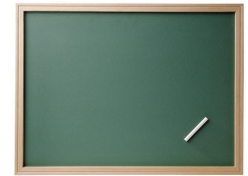
<b>x</b>	<b>c(x)</b>
<b>a</b>	<b>1</b>
<b>b</b>	<b>011</b>

- Are they uniquely decodable?
- Complete prefix code: no unused branches

# WHAT LIMITS DOES UNIQUE DECODABILITY IMPOSE ON THE CODE LENGTH?

- Let's for the moment only think about code lengths (and forget about probabilities)
- ... that is, we'll worry only about unique decodability vs. code length
- If we want our code to be uniquely decodable we cannot allow some codewords to combine to others!
- Short codewords are a scarce resource!

# COST OF CODEWORDS



- Budget of 1
- Cost will be  $1/2^{l(x)}$

0	00	000	0000
			0001
		001	0010
		0011	
	01	010	0100
			0101
011		0110	
	0111		
1	10	100	1000
			1001
		101	1010
		1011	
	11	110	1100
			1101
111		1110	
	1111		

The total symbol code budget

# KRAFT INEQUALITY

- There exists a uniquely decodable code if and only if the codeword lengths satisfy

$$\sum_{x \in \mathcal{A}_X} \frac{1}{2^{l(x)}} \leq 1$$

Furthermore, if such a code exists, a prefix code exists.

- Kraft:** If the inequality is satisfied, then a prefix code exists with the given lengths (and prefix codes are uniquely decodable).
- McMillan:** Conversely, any uniquely decodable code must satisfy this inequality.

# KRAFT INEQUALITY

If equality holds the code is called a “complete code”.

- There exists a uniquely decodable code if and only if the codeword lengths satisfy

$$\sum_{x \in \mathcal{A}_X} \frac{1}{2^{l(x)}} \leq 1$$

Furthermore, if such a code exists, a prefix code exists.

- Kraft:** If the inequality is satisfied, then a prefix code exists with the given lengths (and prefix codes are uniquely decodable).
- McMillan:** Conversely, any uniquely decodable code must satisfy this inequality.

# KRAFT INEQUALITY

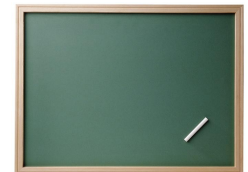
If equality holds the code is called a “complete code”.

- There exists a uniquely decodable code if and only if the codeword lengths satisfy

$$\sum_{x \in \mathcal{A}_X} \frac{1}{2^{l(x)}} \leq 1$$

Furthermore, if such a code exists, a prefix code exists.

- Kraft: If the inequality is satisfied, then a prefix code exists with the given lengths (and prefix codes are uniquely decodable).
- McMillan: Conversely, any uniquely decodable code must satisfy this inequality.



# UNIQUE DECODABILITY VS PREFIX FREE

Kraft inequality holds

$$\sum_{x \in \mathcal{A}_X} \frac{1}{2^{l(x)}} \leq 1$$

# UNIQUE DECODABILITY VS PREFIX FREE

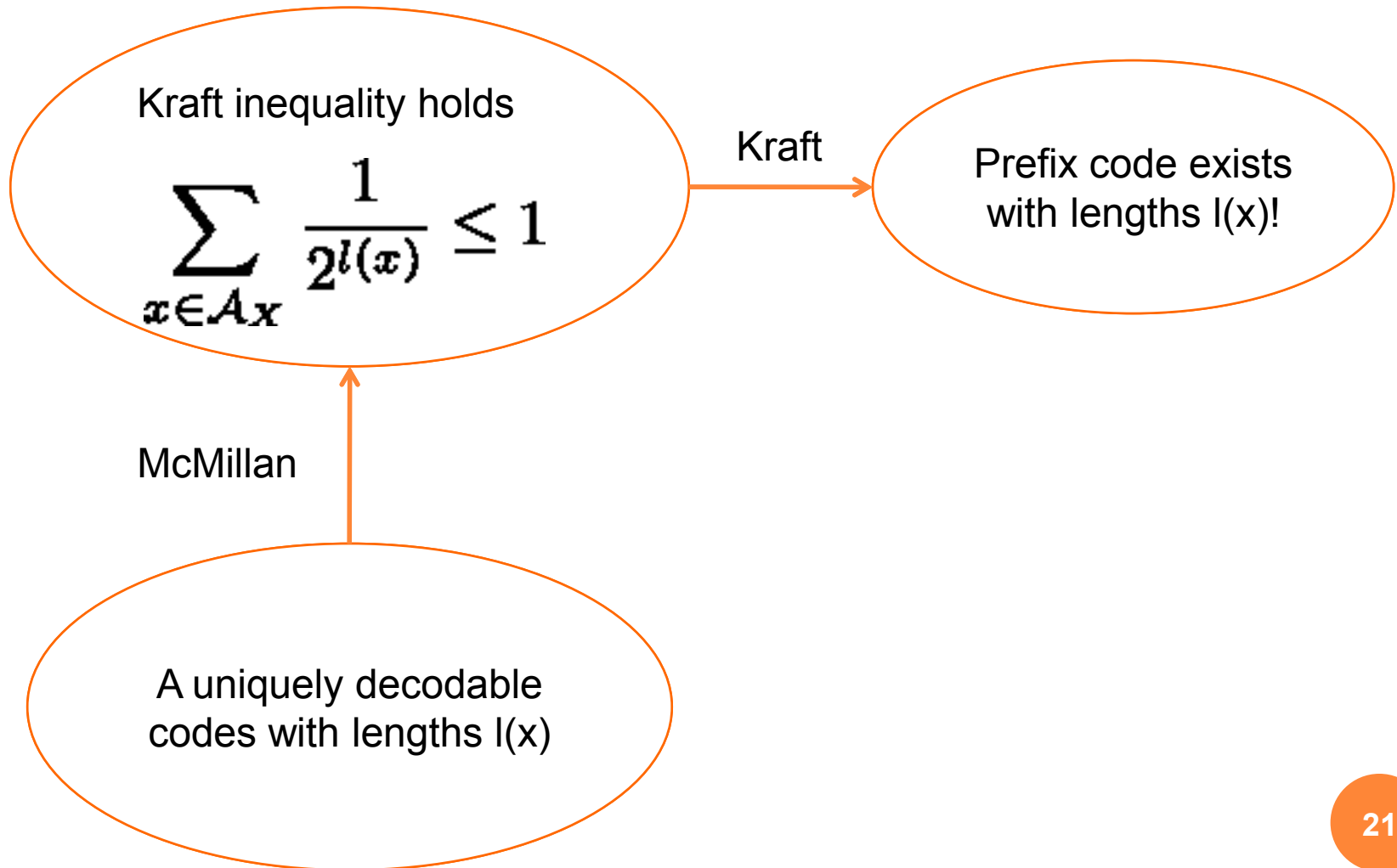
Kraft inequality holds

$$\sum_{x \in \mathcal{A}_x} \frac{1}{2^{l(x)}} \leq 1$$

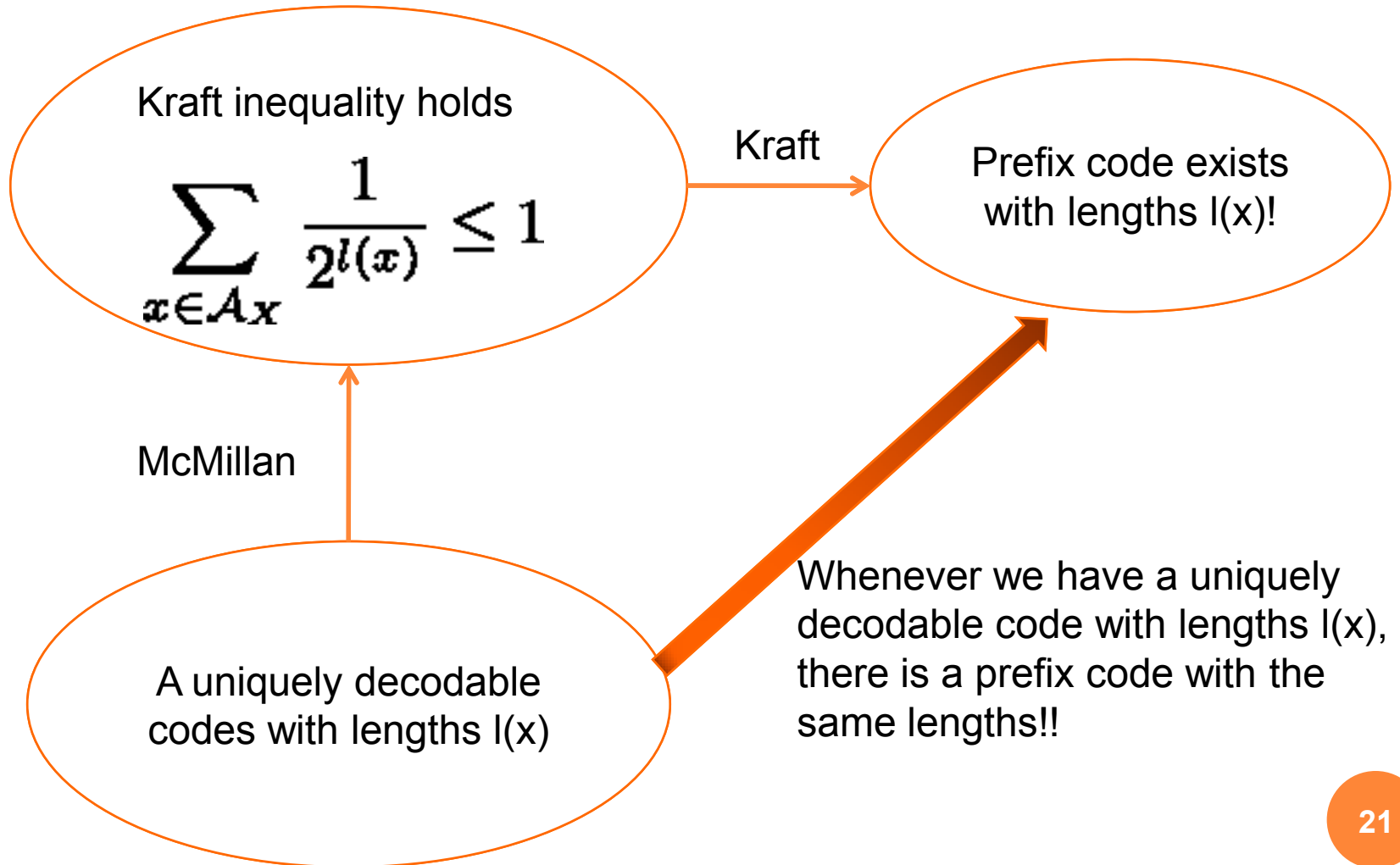
Kraft

Prefix code exists  
with lengths  $l(x)$ !

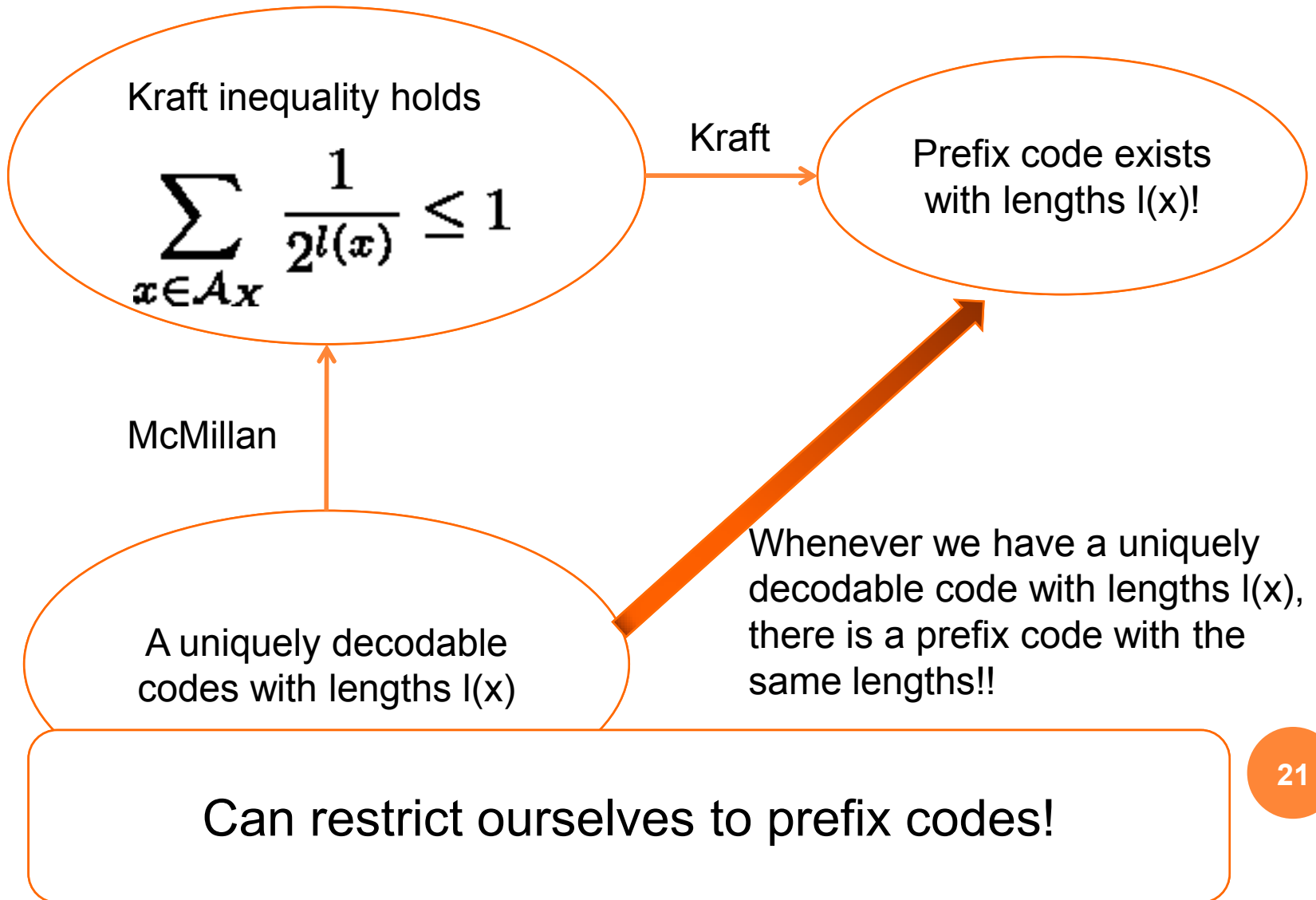
# UNIQUE DECODABILITY VS PREFIX FREE



# UNIQUE DECODABILITY VS PREFIX FREE



# UNIQUE DECODABILITY VS PREFIX FREE



# REQUIREMENTS OF USEFUL SYMBOL CODES

- The extended code must be uniquely decodable
  - Any string must have a unique decoding
- Must be easy to decode
- Achieve the highest possible compression
  - Remember: it's sufficient to consider prefix codes



## EXPECTED CODE LENGTH

- We cannot hope to reduce the total number of codewords, nor the sum of lengths of all codewords – can you say why?
- But we can hope to reduce the expected code length

$$L(C, X) = \sum_{x \in \mathcal{A}_X} P_X(x)l(x)$$

# HOW WELL CAN WE COMPRESS?

- Measured in terms of the average codeword length

$$L(C, X) = \sum_{x \in \mathcal{A}_X} P_X(x) l(x)$$

- What's the best compression we can hope for?
  - Lower bound by the Shannon entropy

$$H(X) \leq L(C, X)$$

# PROOF OF THE LOWER BOUND $H(X) \leq L(C, X)$

- Let's investigate  $L(C, X) = \sum_{x \in \mathcal{A}_X} P_X(x) l(x)$

# PROOF OF THE LOWER BOUND $H(X) \leq L(C, X)$

- Let's investigate  $L(C, X) = \sum_{x \in \mathcal{A}_X} P_X(x) l(x)$

Actual distribution

$$P(x)$$

# PROOF OF THE LOWER BOUND $H(X) \leq L(C, X)$

- Let's investigate  $L(C, X) = \sum_{x \in \mathcal{A}_X} P_X(x) l(x)$

Actual distribution

$$P(x)$$

Implicit distribution

$$Q(x) = \frac{2^{-\ell(x)}}{z}$$
$$z = \sum_x 2^{-\ell(x)}$$

# PROOF OF THE LOWER BOUND $H(X) \leq L(C, X)$

- Let's investigate  $L(C, X) = \sum_{x \in \mathcal{A}_X} P_X(x) l(x)$

Actual distribution

$$P(x)$$

Implicit distribution

$$Q(x) = \frac{2^{-\ell(x)}}{z}$$
$$z = \sum_x 2^{-\ell(x)}$$

Difference measured by the relative entropy “distance”  
(Kullback-Leibler divergence)

$$D(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

# PROOF OF THE LOWER BOUND $H(X) \leq L(C, X)$

- Let's investigate  $L(C, X) = \sum_{x \in \mathcal{A}_X} P_X(x) l(x)$

Actual distribution

$$P(x)$$

Implicit distribution

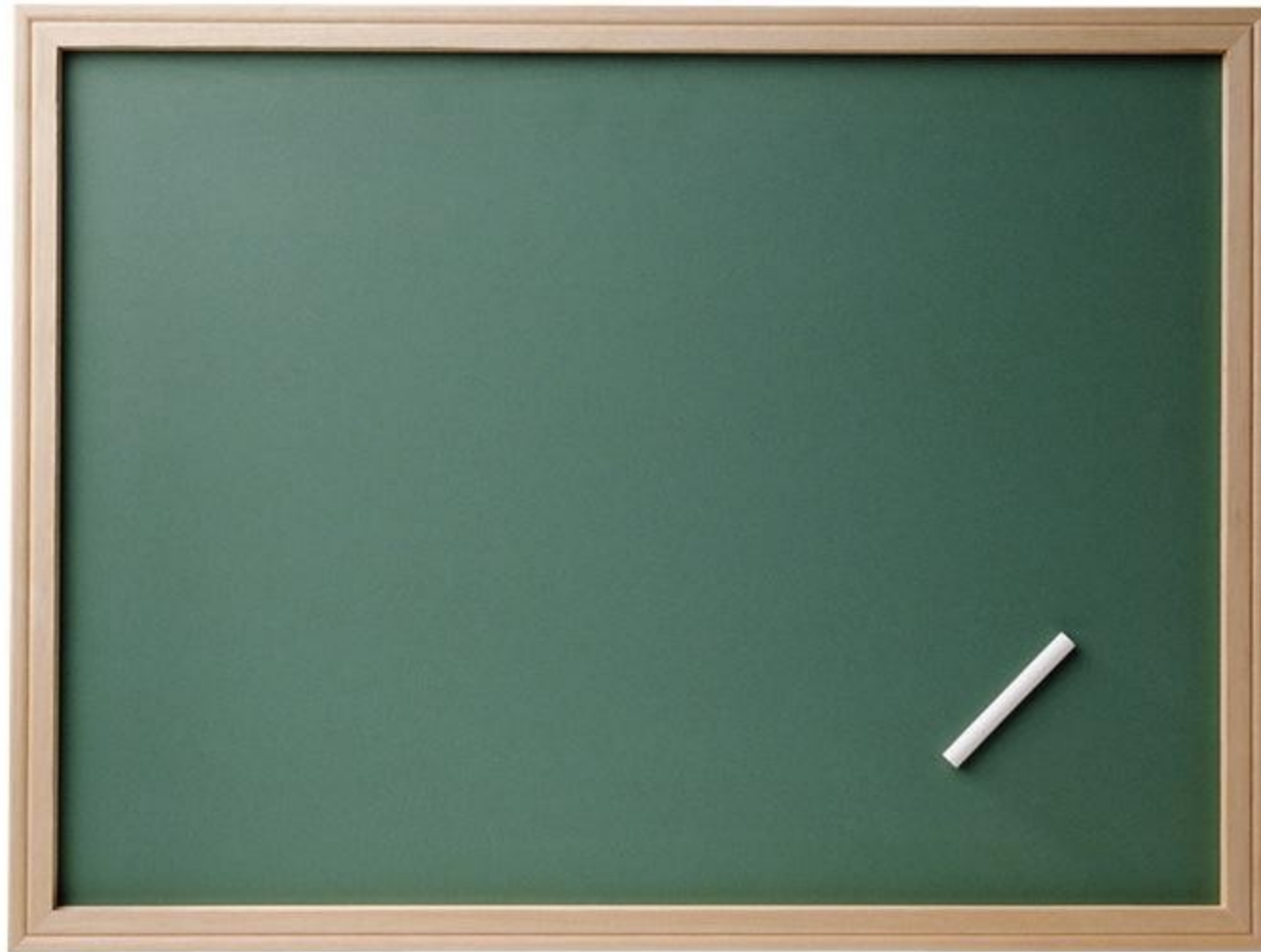
$$Q(x) = \frac{2^{-\ell(x)}}{z}$$
$$z = \sum_x 2^{-\ell(x)}$$

Difference measured by the relative entropy “distance”  
(Kullback-Leibler divergence)

$$D(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} \geq 0$$

Gibbs inequality

# PROOF OF THE LOWER BOUND $H(X) \leq L(C, X)$



## SUMMARY OF REMARKS ON THE LOWER BOUND

- When is the lower bound achieved?

$$H(X) = L(C, X)$$

- Optimal source codelength equal to information content

$$l(x) = \log \frac{1}{P_X(x)}$$

- Complete code  $z = 1$

- Conversely, any given codelengths would be optimal for the distribution

$$Q_X(x) = \frac{2^{-l(x)}}{z}$$

$$z = \sum_{\mathbf{x}} 2^{-l(\mathbf{x})}$$

WHAT IF WE USE A COMPLETE CODE BUT  
DIFFERENT CODELENGTHS?

# WHAT IF WE USE A COMPLETE CODE BUT DIFFERENT CODELENGTHS?

Average codeword length becomes larger

$$L(C, X) = H(X) + \sum_x P_X(x) \log \frac{P_X(x)}{\tilde{Q}_X(x)} \quad \tilde{Q}_X(x) = 2^{-l(x)}$$

# WHAT IF WE USE A COMPLETE CODE BUT DIFFERENT CODELENGTHS?

Average codeword length becomes larger

$$L(C, X) = H(X) + \sum_x P_X(x) \log \frac{P_X(x)}{\tilde{Q}_X(x)} \quad \tilde{Q}_X(x) = 2^{-l(x)}$$



Term is positive and equal to the so-called relative entropy

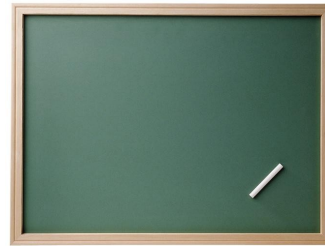
$$D(P||\tilde{Q})$$

# HOW MUCH CAN WE COMPRESS?

- There exists a compression scheme such that

$$L(C, X) \leq H(X) + 1$$

- Proof using the Kraft inequality



# WHAT WE'LL DO TODAY

- Lossless data compression
  - Basic concepts
  - Symbol codes
  - Kraft's inequality
  - Lossless compression and entropy
  - Shannon's source coding theorem for symbol codes
- An optimal symbol code: Huffman coding
- Limits of symbol codes

# GOAL

- We want an explicit coding scheme such that:

$$H(X) \leq L(C, X) < H(X) + 1$$

- Ideally, *no other symbol code has lower expected codelength*
- We call such a scheme *optimal*
- Any ideas?

# SOME CLUES

- We know that we only need to think about prefix codes
- ... and these will automatically be uniquely decodable
- Any prefix code is a tree.
- So how can we make such a tree?

# ATTEMPT #1: TOP DOWN (GREEDY)

- Example

$$\mathcal{A}_X = \{a, b, c, d, e, f, g\}$$

$$P_X = \{0.01, 0.24, 0.05, 0.20, 0.47, 0.01, 0.02\}$$

$$H(X) = 1.93$$

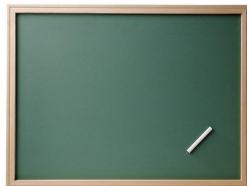
# ATTEMPT #1: TOP DOWN (GREEDY)

- Example

$$\mathcal{A}_X = \{a, b, c, d, e, f, g\}$$

$$P_X = \{0.01, 0.24, 0.05, 0.20, 0.47, 0.01, 0.02\}$$

$$H(X) = 1.93$$



$$L(C, X) = 2.53$$

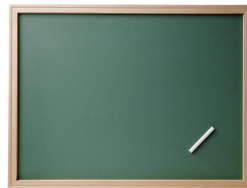
## ATTEMPT #1: TOP DOWN (GREEDY)

- Example

$$\mathcal{A}_X = \{a, b, c, d, e, f, g\}$$

$$P_X = \{0.01, 0.24, 0.05, 0.20, 0.47, 0.01, 0.02\}$$

$$H(X) = 1.93$$



# ATTEMPT #2: BOTTOM UP (HUFFMAN)

## ATTEMPT #2: BOTTOM UP (HUFFMAN)

- Take the two least probable symbols
  - These will get the longest codewords
    - Of equal length, differing in the last digit only

## ATTEMPT #2: BOTTOM UP (HUFFMAN)

- Take the two least probable symbols
  - These will get the longest codewords
    - Of equal length, differing in the last digit only
- Combine them into a single symbol

## ATTEMPT #2: BOTTOM UP (HUFFMAN)

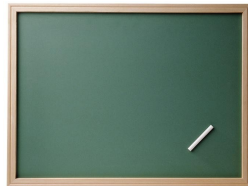
- Take the two least probable symbols
  - These will get the longest codewords
    - Of equal length, differing in the last digit only
- Combine them into a single symbol

Repeat

## ATTEMPT #2: BOTTOM UP (HUFFMAN)

- Take the two least probable symbols
  - These will get the longest codewords
    - Of equal length, differing in the last digit only
- Combine them into a single symbol

Repeat

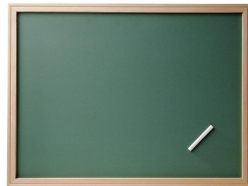


$$L(C, X) = 1.97$$

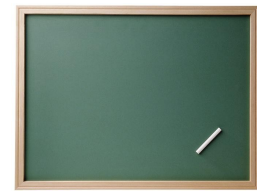
## ATTEMPT #2: BOTTOM UP (HUFFMAN)

- Take the two least probable symbols
  - These will get the longest codewords
    - Of equal length, differing in the last digit only
- Combine them into a single symbol

Repeat

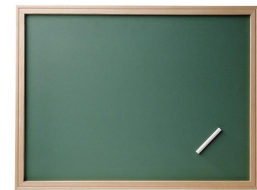


# HUFFMAN CODING EXAMPLE



$\mathcal{A}_X$	$a$	$b$	$c$	$d$	$e$
$P_X$	0.25	0.25	0.2	0.15	0.15

# HUFFMAN CODING EXAMPLE



$\mathcal{A}_X$	$a$	$b$	$c$	$d$	$e$
$P_X$	0.25	0.25	0.2	0.15	0.15

$x$	$P_X(x)$	$h(p_X(x))$	$l(x)$	$c(x)$
$a$	0.25	2	2	00
$b$	0.25	2	2	10
$c$	0.2	2.3	2	11
$d$	0.15	2.7	3	010
$e$	0.15	2.7	3	011

# PERFORMANCE OF HUFFMAN CODING

- Can be proven to be optimal! No other better symbol code
- Does this mean it is the solution to all our compression problems?

# DISADVANTAGES

- Does not make use of context!
  - E.g., when compressing English language

# DISADVANTAGES

- Does not make use of context!
  - E.g., when compressing English language

That is, does not adapt to changing ensemble

# DISADVANTAGES

- Does not make use of context!
  - E.g., when compressing English language

That is, does not adapt to changing ensemble

- Single bit can be significant when the entropy is small

$$H(X) \leq L(C, X) < H(X) + 1$$

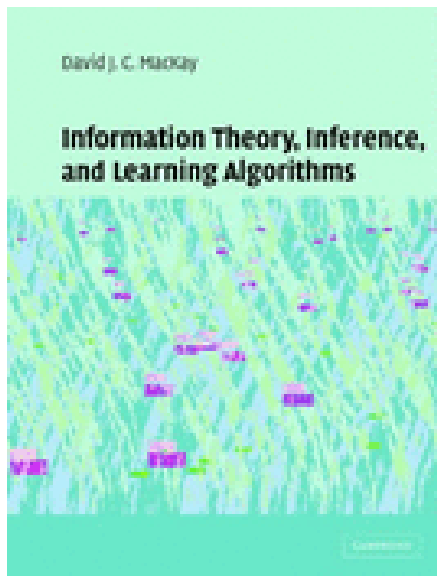
Can be dealt with to some extent by compressing blocks at once

# LET'S SUMMARIZE

- Basics of lossless compression and symbol codes
- Prefix codes
- Kraft's inequality
- Optimal expected code length determined by the Shannon entropy
- Optimal Huffman code
- Limits of symbol codes: no context, at least one bit per encoded symbol
  
- Next time: Stream codes

# READING FOR THIS LECTURE

- Chapter 5 in the book



Information Theory, Inference and  
Learning Algorithms  
by David J. C. MacKay  
Cambridge University Press, 2003

- Homework due by Monday 2pm next week, **upload to IVLE Workbin**