

Algorithms

Lecture 4: Greedy Algorithm

Anxiao (Andrew) Jiang

CH 16. Greedy Algorithms

16.3 Huffman Code

Symbol	a	b	c	d	e	f
Probability	0.45	0.13	0.12	0.16	0.09	0.05

How to represent the symbols using bits, to minimize the average number of bits needed?

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05

How to represent the symbols using bits, to minimize the average number of bits needed?

Assume: we use “Fixed Length Code (FLC)”

1-bit codewords: 0, 1 (not enough)

2-bit codewords: 00, 01, 10, 11 (not enough)

3-bit codewords: 000, 001, 010, 011, 100, 101, 110, 111 (enough)

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

How to represent the symbols using bits, to minimize the average number of bits needed?

Assume: we use “Fixed Length Code (FLC)”

1-bit codewords: 0, 1 (not enough)

2-bit codewords: 00, 01, 10, 11 (not enough)

3-bit codewords: 000, 001, 010, 011, 100, 101, 110, 111 (enough)

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Compression (encoding): turn a text to bits

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a* *a b f e d c a ...*

Compression (encoding): turn a text to bits

000

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000000

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101 100

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101 100 011

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101 100 011 010

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a* ...

Compression (encoding): turn a text to bits

000 000 001 101 100 011 010 000 ...

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101 100 011 010 000 ...

Decompression (decoding): turn a bit sequence back to text

a

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101 100 011 010 000 ...

Decompression (decoding): turn a bit sequence back to text

a a

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101 100 011 010 000 ...

Decompression (decoding): turn a bit sequence back to text

a a b

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101 100 011 010 000 ...

Decompression (decoding): turn a bit sequence back to text

a a b f

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101 100 011 010 000 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101 100 011 010 000 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e d

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101 100 011 010 000 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e d c

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101 100 011 010 000 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e d c a ...

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

000 000 001 101 100 011 010 000 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e d c a ...

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101

3 bits per symbol

Variable Length Code (VLC): the codewords can have different lengths.

Prefix Code: no codeword is the prefix of another codeword.

Let's study Variable Length Prefix Code (VLPC).

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Variable Length Code (VLC): the codewords can have different lengths.

Prefix Code: no codeword is the prefix of another codeword.

Let's study Variable Length Prefix Code (VLPC).

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

00

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a* ...

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a *a*

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 **1** 1 0 0 1 1 0 1 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b f

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b f

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b f

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b f

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e d

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e d

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e d

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e d c

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
FLC	000	001	010	011	100	101
VLPC	0	101	100	111	1101	1100

3 bits per symbol

Text: *a a b f e d c a ...*

Compression (encoding): turn a text to bits

0 0 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0 0 ...

Decompression (decoding): turn a bit sequence back to text

a a b f e d c a ...

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	
Probability	0.45	0.13	0.12	0.16	0.09	0.05	
FLC	000	001	010	011	100	101	3 bits per symbol
VLPC	0	101	100	111	1101	1100	2.24 bits per symbol

Average codeword length for VLPC:

$$1 \times 0.45 + 3 \times 0.13 + 3 \times 0.12 + 3 \times 0.16 + 4 \times 0.09 + 4 \times 0.05 = 2.24 \text{ bits/symbol}$$

16.3 Huffman Code

Input: n symbols s_1, s_2, \dots, s_n .

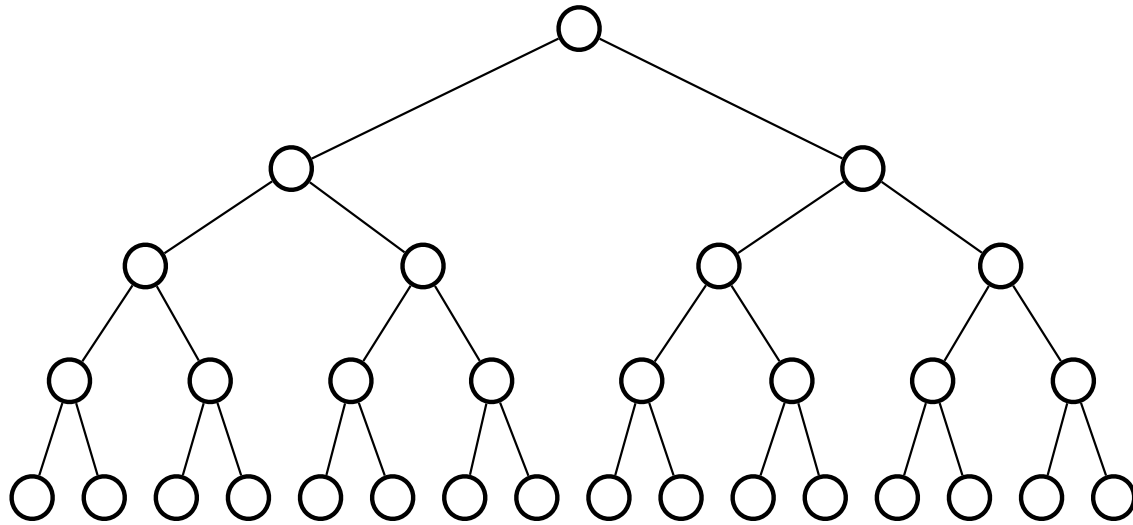
For $i = 1, 2, \dots, n$, the symbol s_i has probability f_i .

Output: Design a **prefix code** for the n symbols such that the **average codeword length** is minimized.

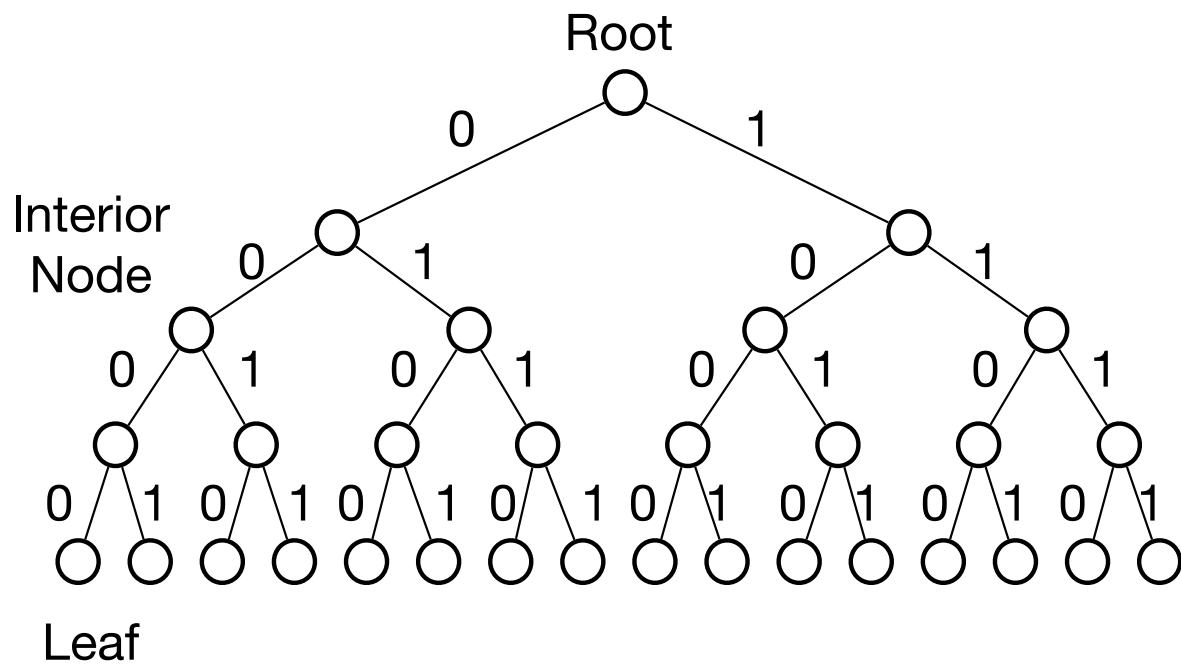
Assume the codeword for symbol s_i has L_i bits.

$$\text{Average Codeword Length} = \sum_{i=1}^n f_i L_i$$

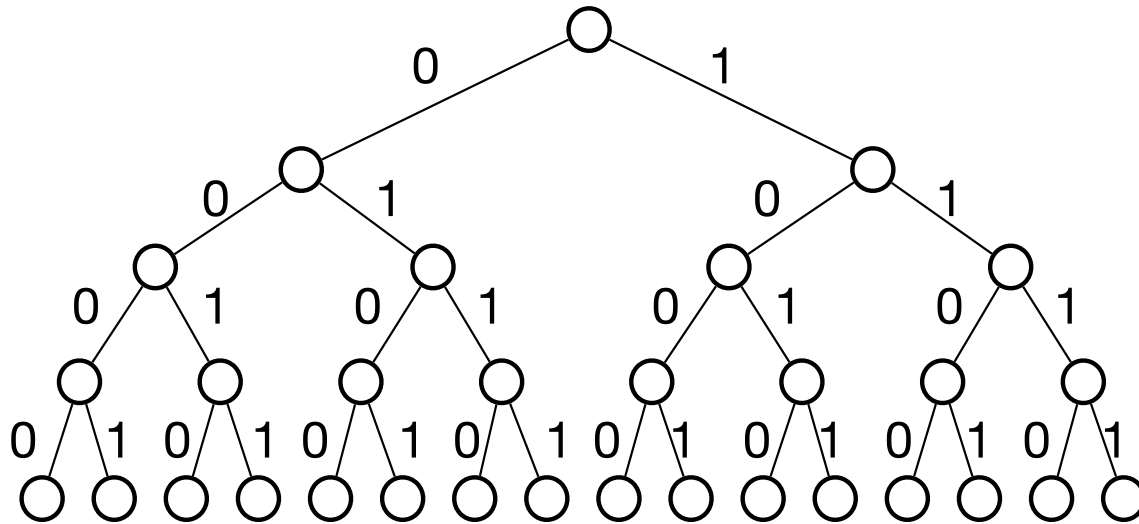
16.3 Huffman Code



16.3 Huffman Code

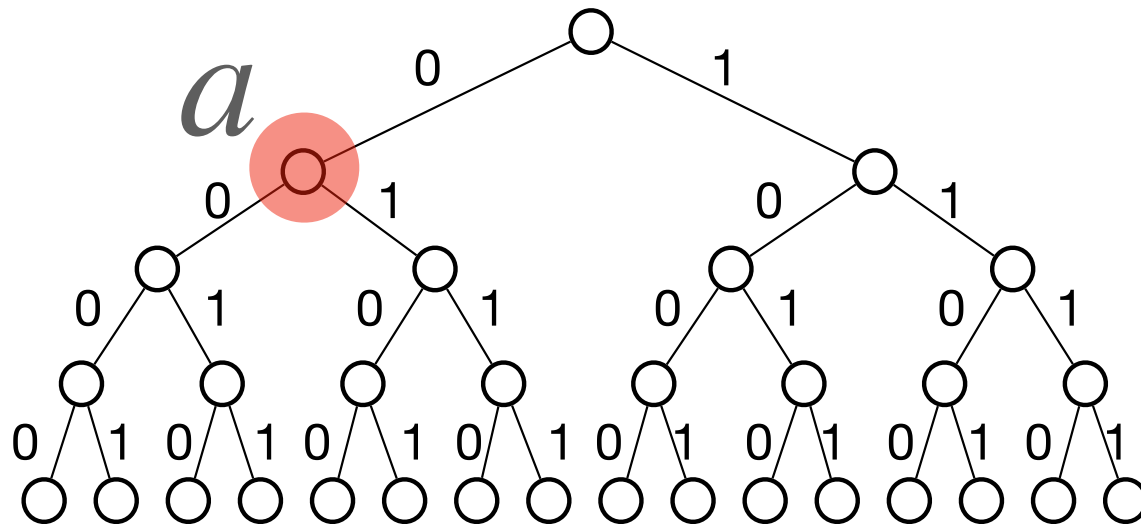


16.3 Huffman Code



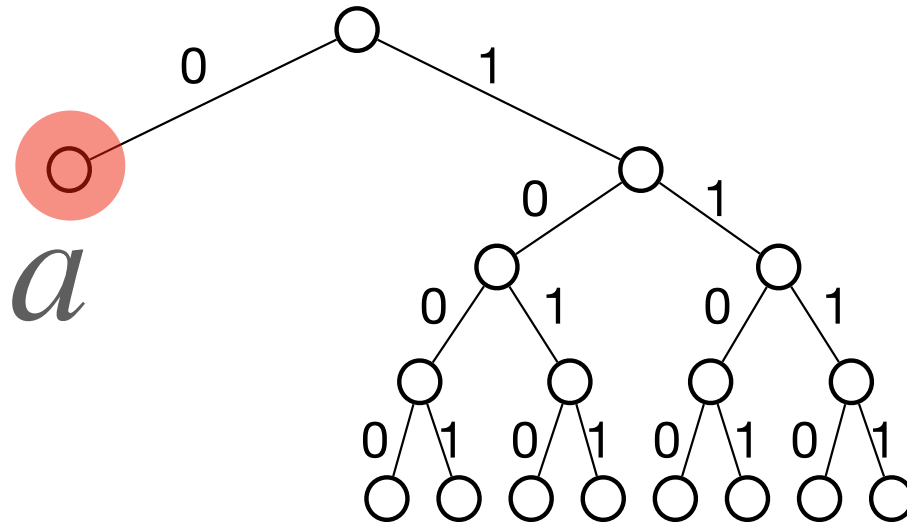
Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



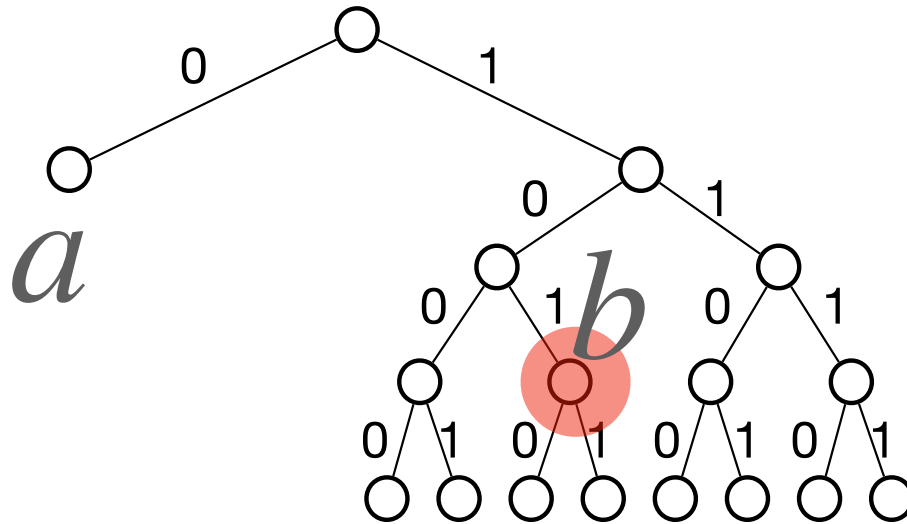
Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



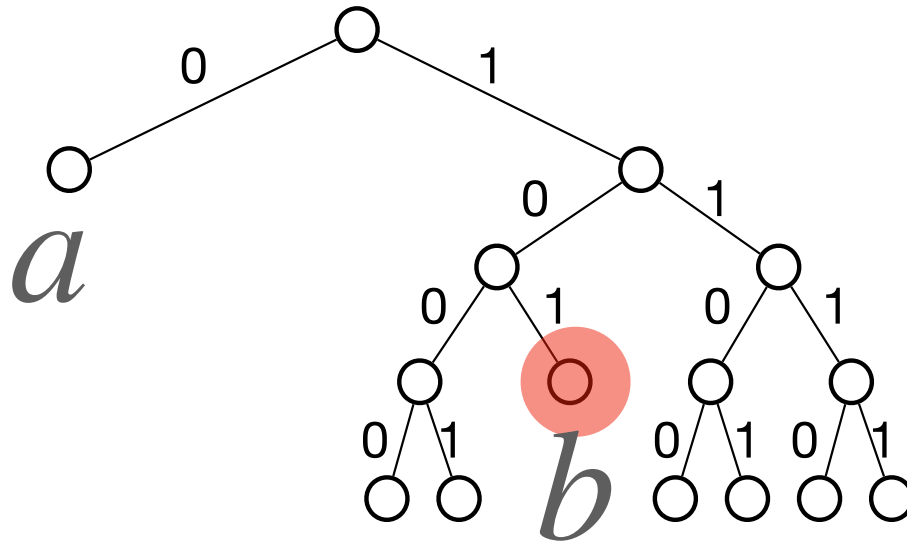
Symbol	a	b	c	d	e	f
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



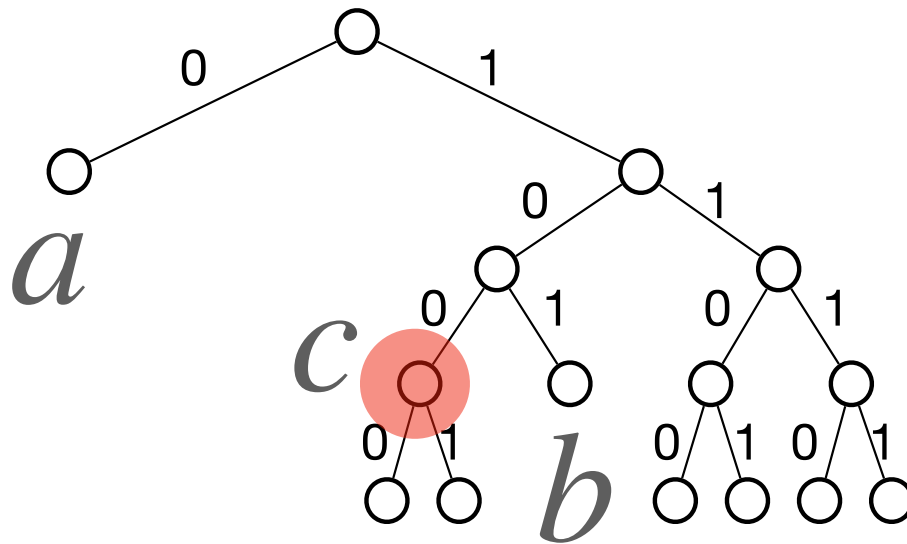
Symbol	a	b	c	d	e	f
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



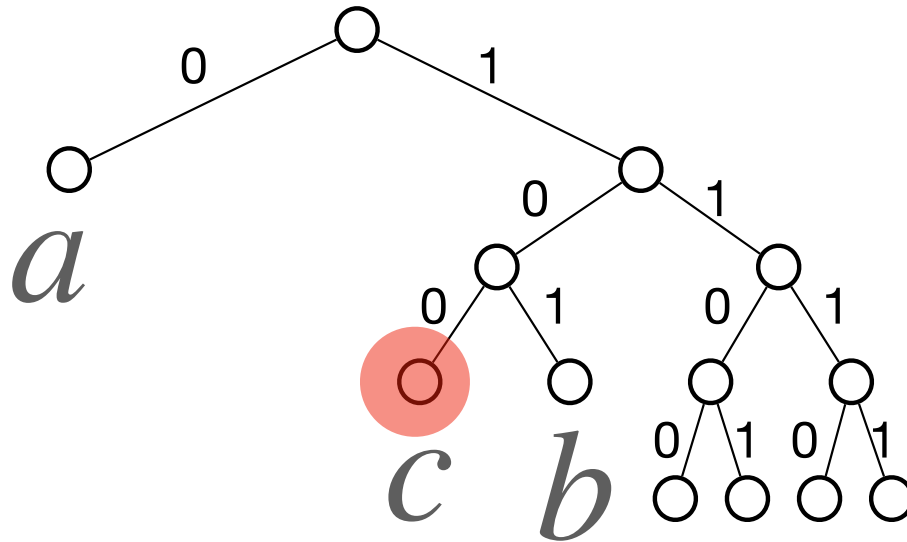
Symbol	a	b	c	d	e	f
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



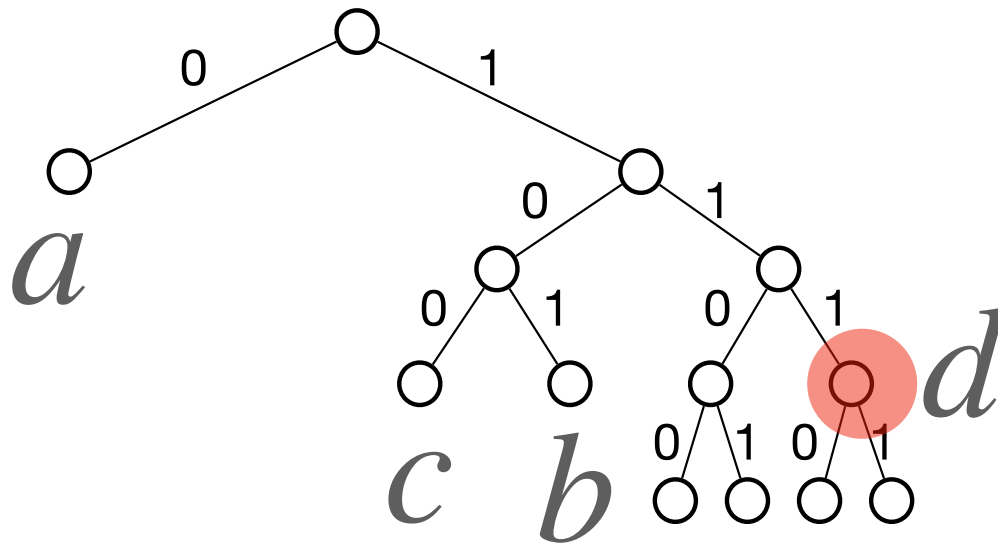
Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



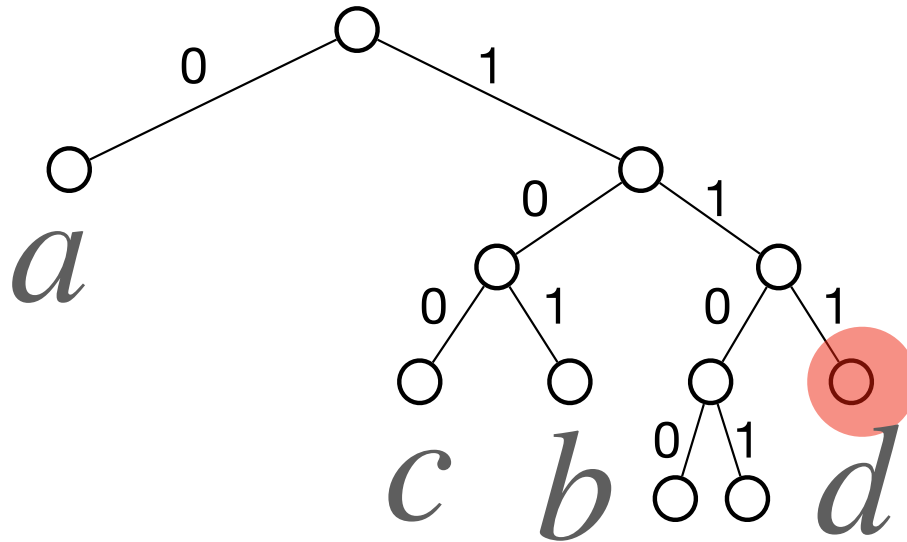
Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



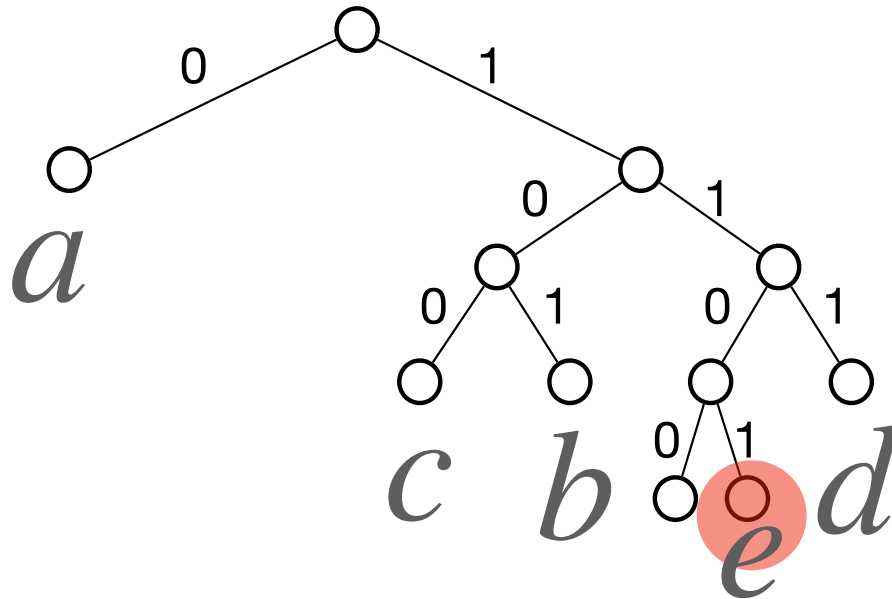
Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



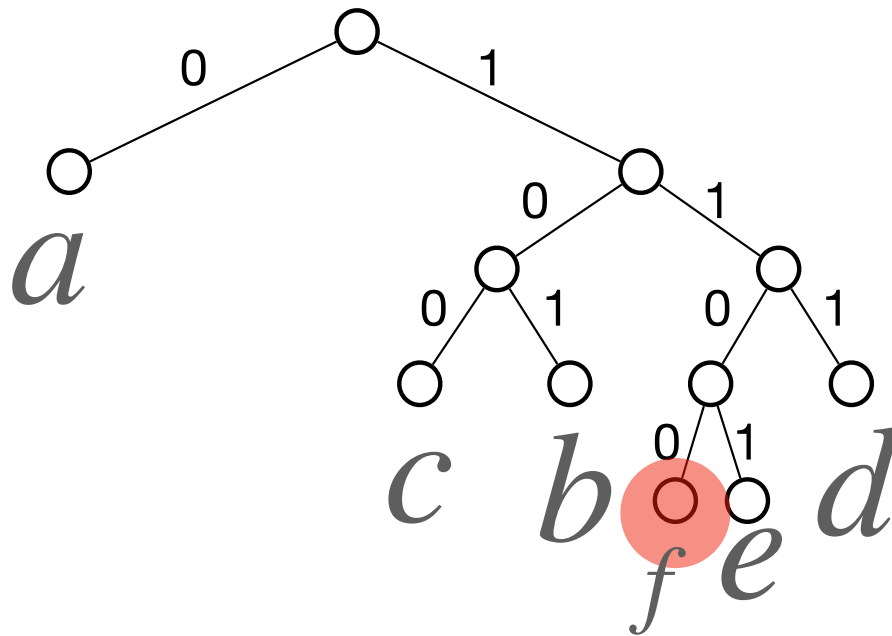
Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



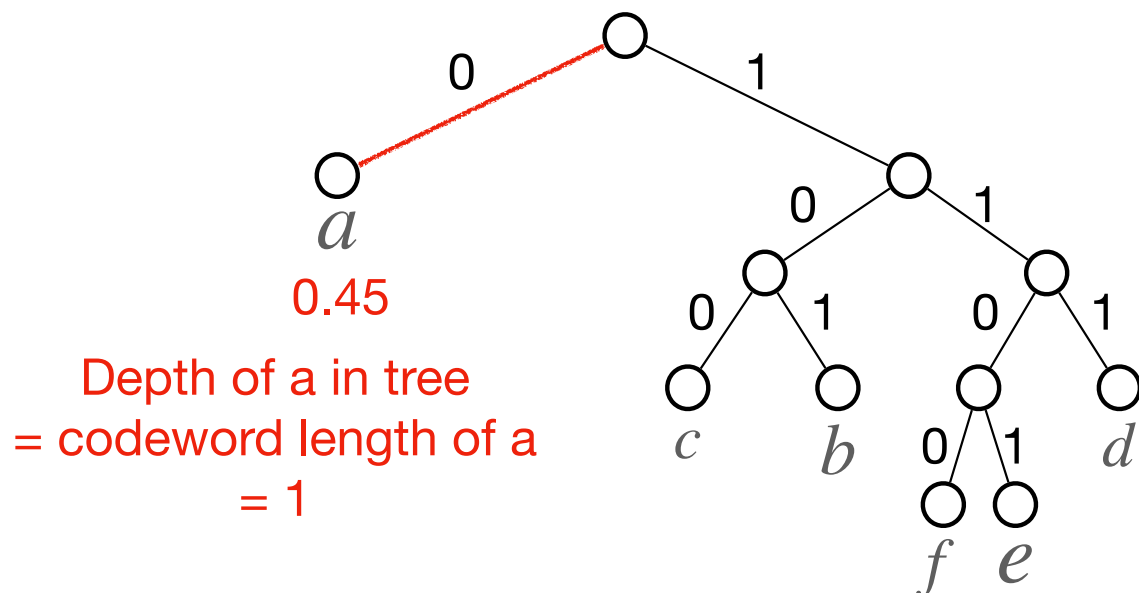
Leaves of a subtree



Prefix code

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



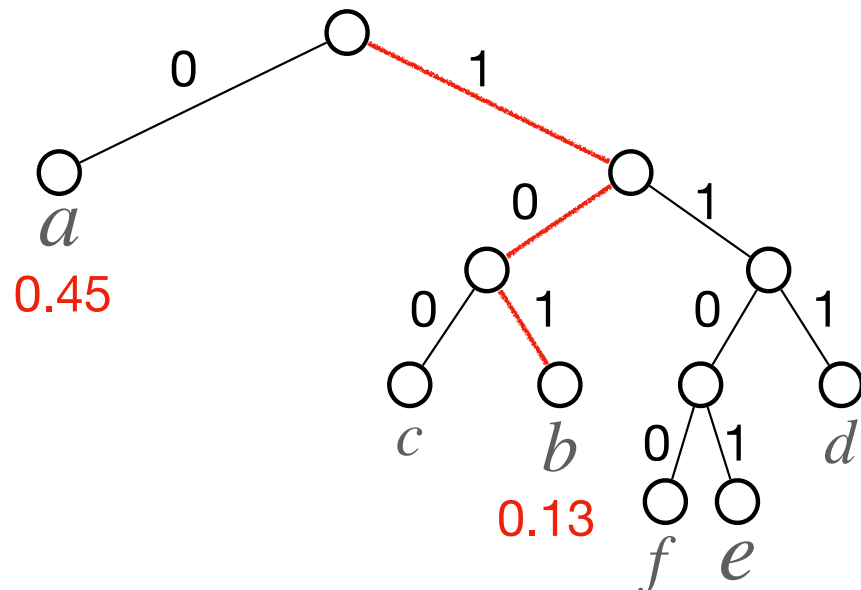
Average **codeword length**



Average **depth of leaves**

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



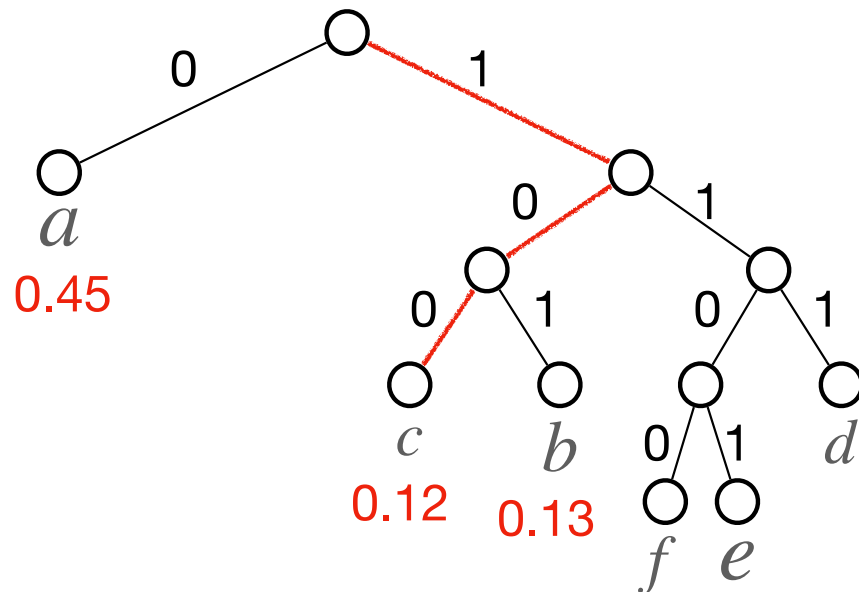
Average **codeword length**



Average **depth of leaves**

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



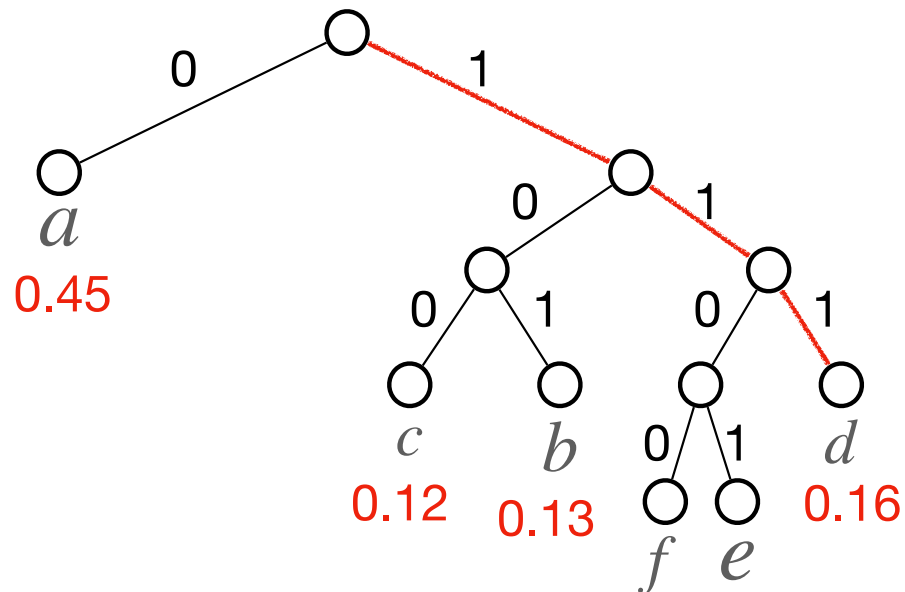
Average **codeword length**



Average **depth of leaves**

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



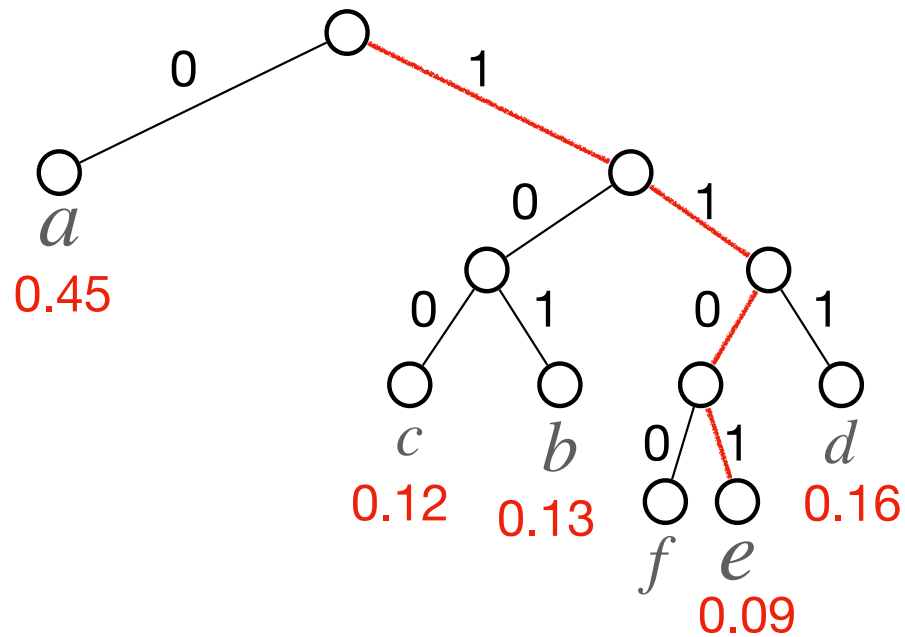
Average **codeword length**



Average **depth of leaves**

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



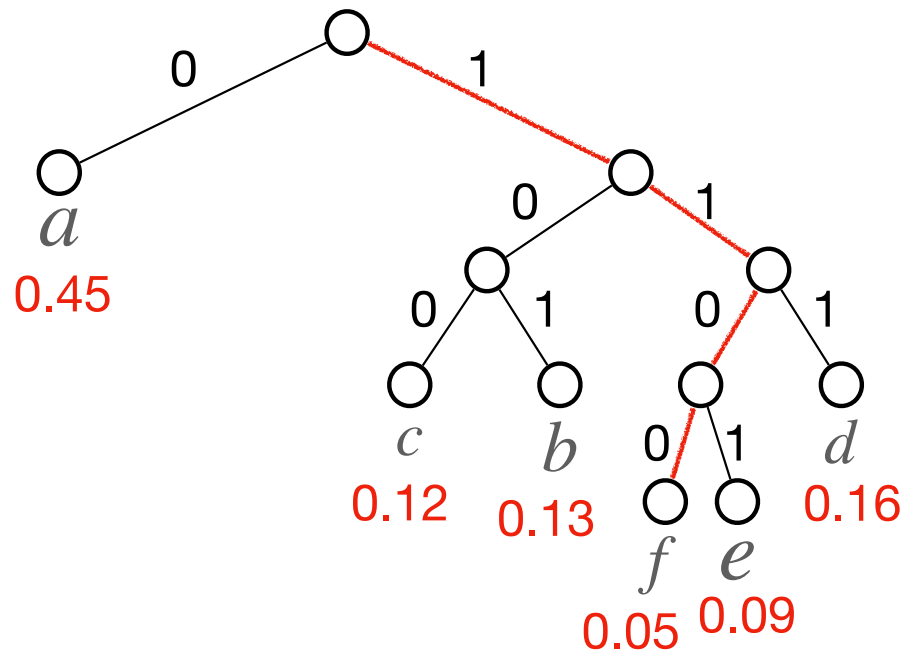
Average **codeword length**



Average **depth of leaves**

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



Average **codeword length**

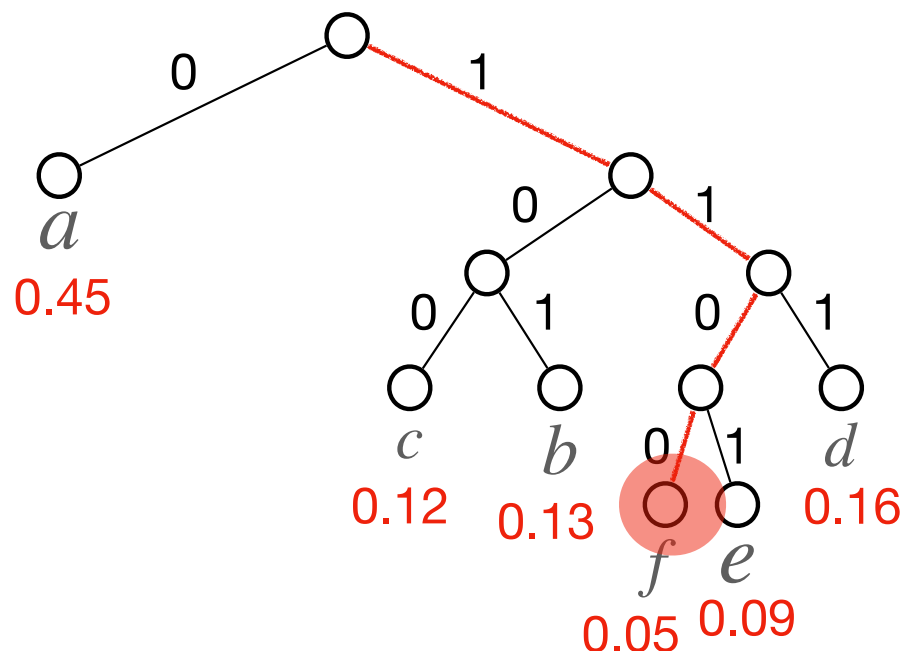


Average **depth of leaves**

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

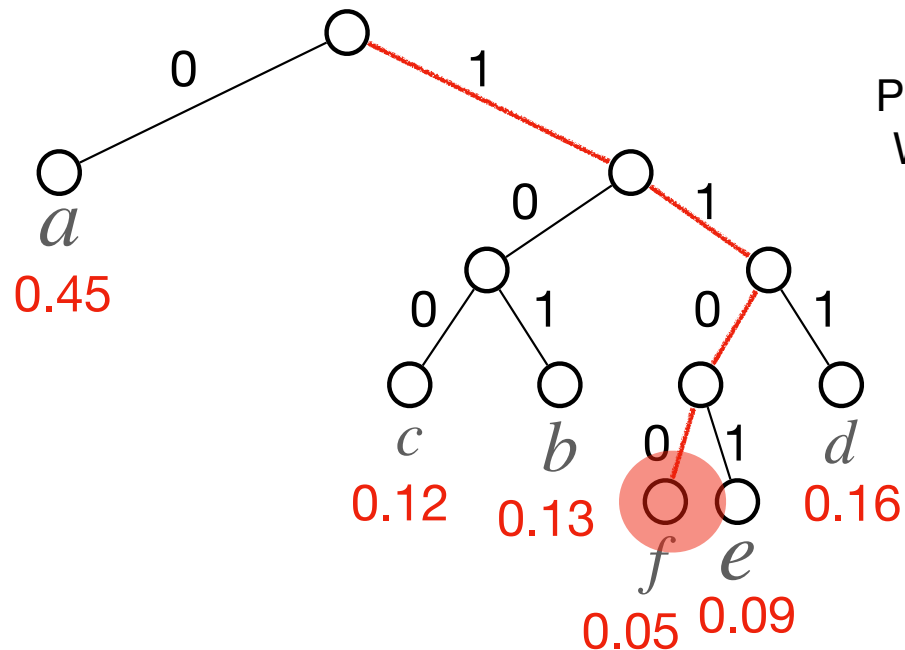
16.3 Huffman Code

Property 1 of optimal code:
The (or one) symbol of **lowest probability**
has the **longest codeword**.



Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code

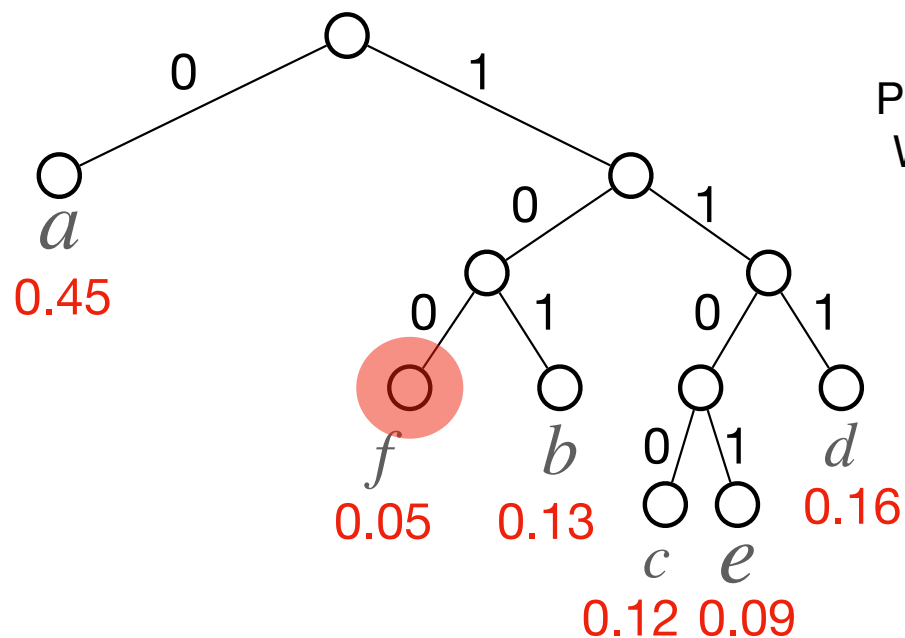


Property 1 of optimal code:
The (or one) symbol of **lowest probability**
has the **longest codeword**.

Proof: If not, then we can switch its codeword
With another symbol, and get a better code.
That will be a contradiction.

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code

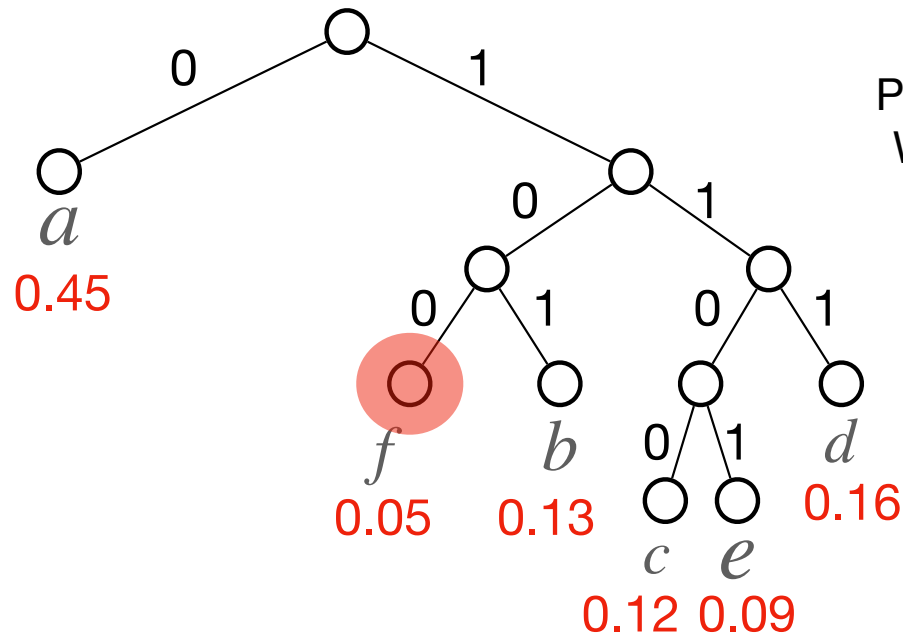


Property 1 of optimal code:
The (or one) symbol of **lowest probability**
has the **longest codeword**.

Proof: If not, then we can switch its codeword
With another symbol, and get a better code.
That will be a contradiction.

Example: Assume *f* does not have
the longest codeword.

16.3 Huffman Code



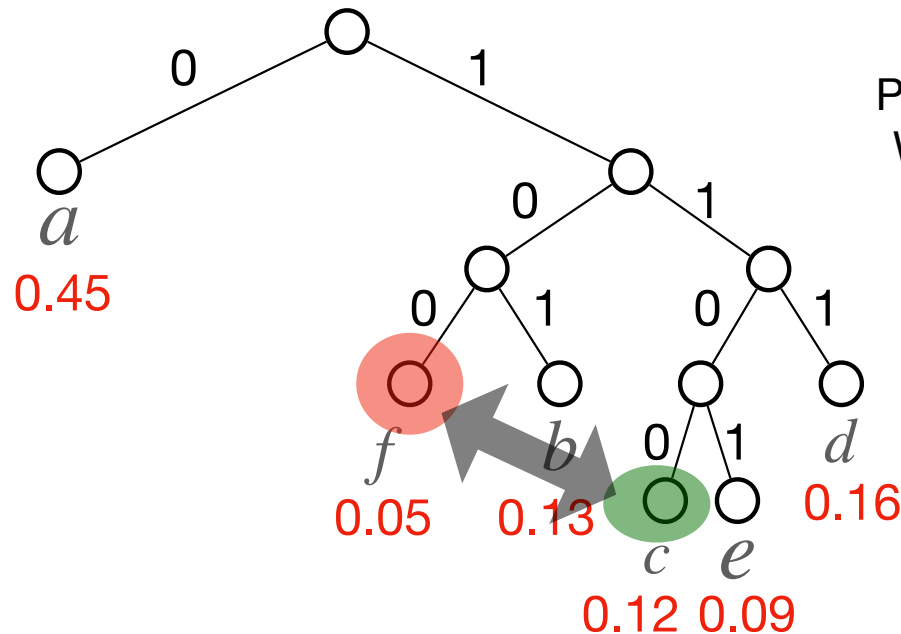
Property 1 of optimal code:
The (or one) symbol of **lowest probability**
has the **longest codeword**.

Proof: If not, then we can switch its codeword
With another symbol, and get a better code.
That will be a contradiction.

Example: Assume *f* does not have
the longest codeword.

Average codeword length = $1 \times 0.45 + 3 \times 0.13 + 4 \times 0.12 + 3 \times 0.16 + 4 \times 0.09 + 3 \times 0.05$

16.3 Huffman Code



Property 1 of optimal code:
The (or one) symbol of **lowest probability**
has the **longest codeword**.

Proof: If not, then we can switch its codeword
With another symbol, and get a better code.
That will be a contradiction.

Example: Assume *f* does not have
the longest codeword.

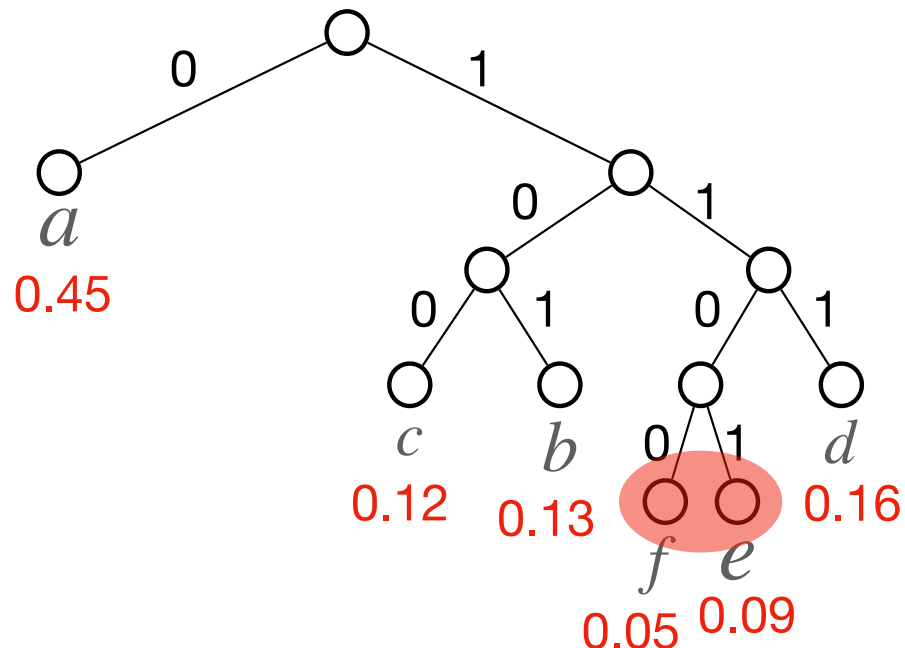
Now switch *a* with *c*.

Average codeword length = $1 \times 0.45 + 3 \times 0.13 + 4 \times 0.12 + 3 \times 0.16 + 4 \times 0.09 + 3 \times 0.05$

After switch, average codeword length = $1 \times 0.45 + 3 \times 0.13 + 3 \times 0.12 + 3 \times 0.16 + 4 \times 0.09 + 4 \times 0.05$

$(4 \times 0.12 + 3 \times 0.05) - (3 \times 0.12 + 4 \times 0.05) = (4 - 3)(0.12 - 0.05) > 0$ **Switch makes code better!**

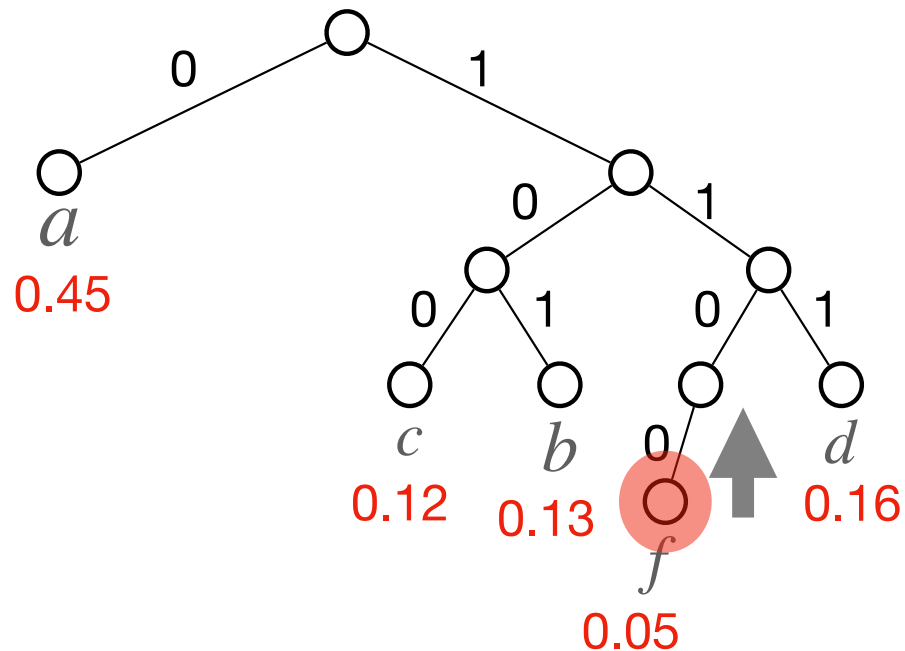
16.3 Huffman Code



Property 2 of optimal code:
The symbol of lowest probability
and longest codeword
has a **sibling leaf node**.

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code

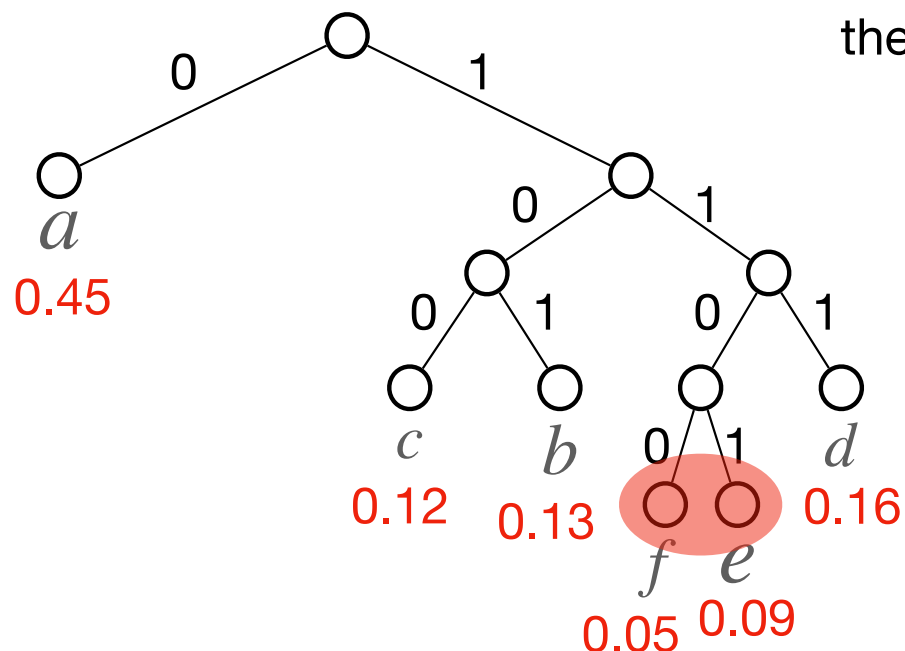


Property 2 of optimal code:
The symbol of lowest probability
and longest codeword
has a **sibling leaf node**.

Proof: If not, we can move the
codeword up to make it shorter.

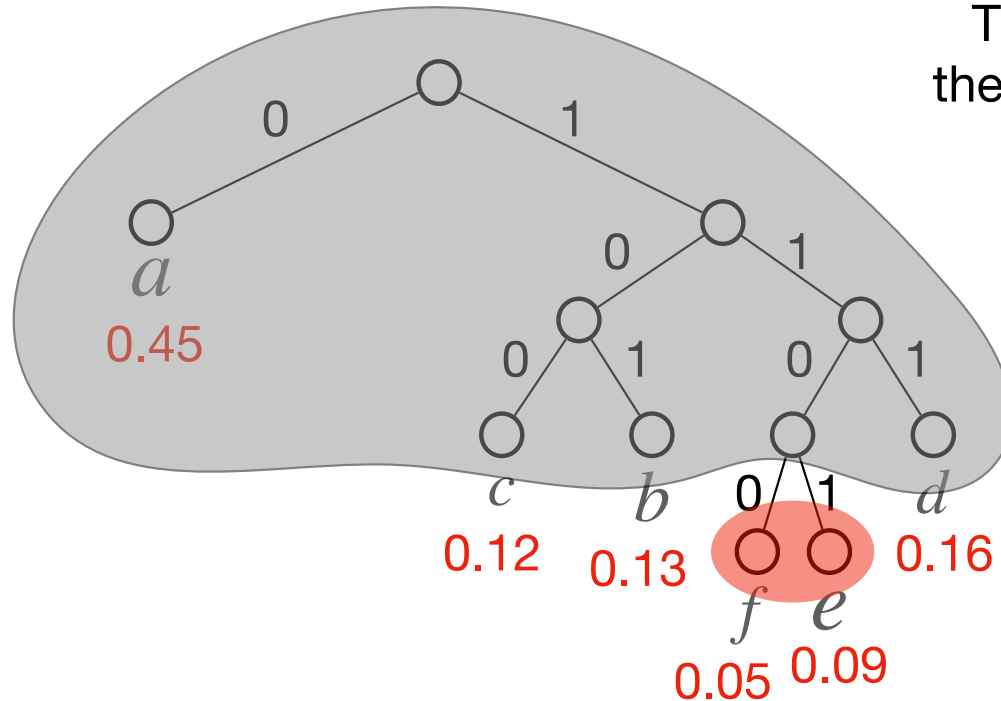
16.3 Huffman Code

Property 3 of optimal code:
There exists an optimal code where
the **two symbols of lowest probabilities**
are siblings.



Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05
VLPC	0	101	100	111	1101	1100

16.3 Huffman Code



Property 3 of optimal code:

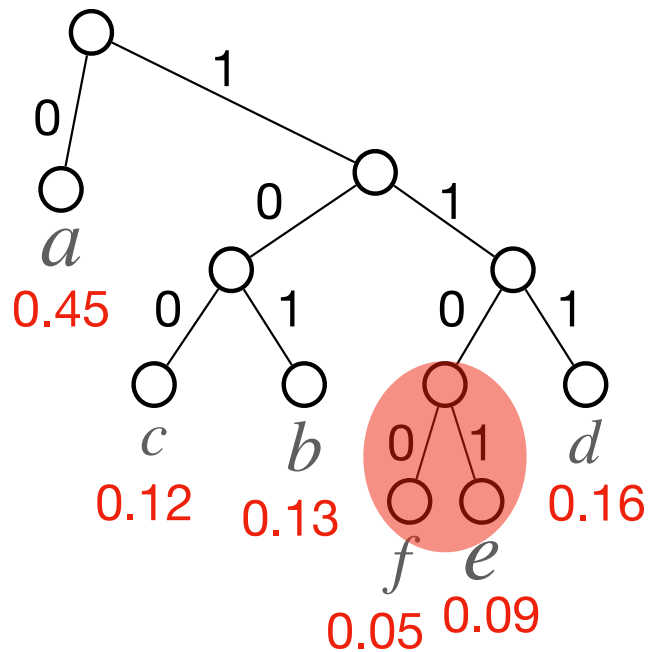
There exists an optimal code where the **two symbols of lowest probabilities** are siblings.

So we can build an optimal tree this way: first put the two symbols of lowest probabilities as siblings. Then figure out the rest of the tree.

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05

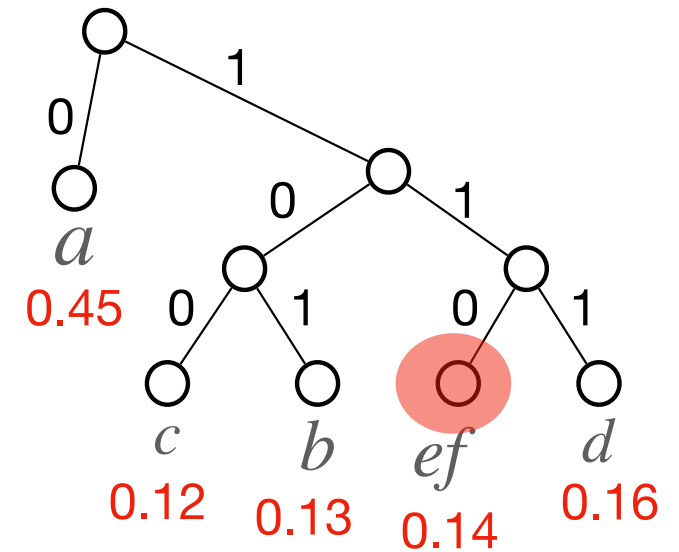
Idea: Once we put the two symbols are siblings,
see them as one symbol (node) and combine their probabilities.

Original Huffman Code



What is
the relationship
between
the two
Huffman codes?

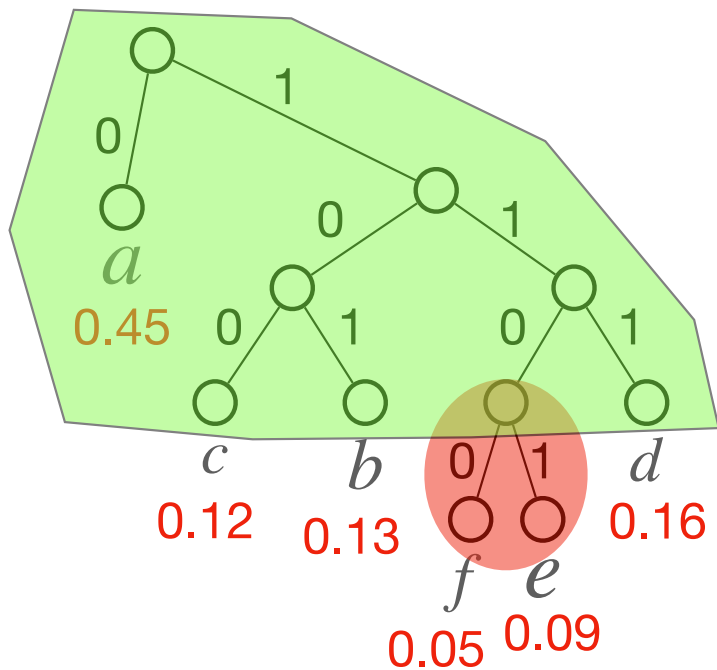
New Huffman Code



Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>ef</i>
Probability	0.45	0.13	0.12	0.16	0.14

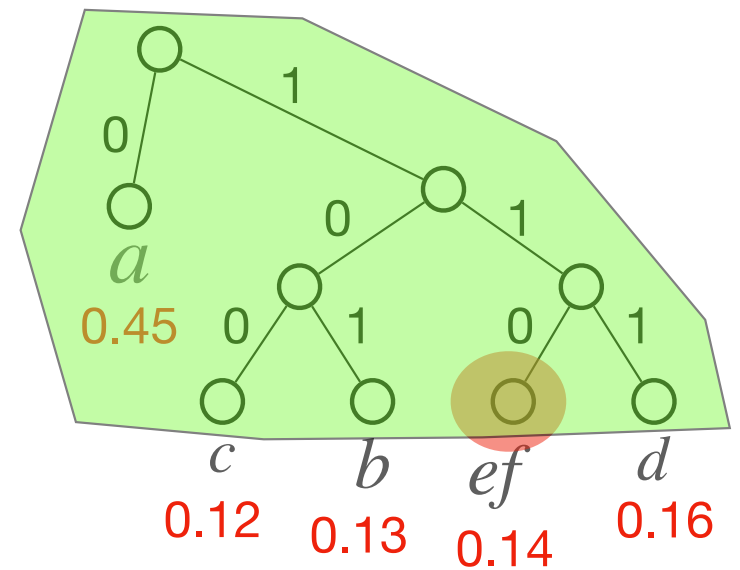
Original Code



What is
the relationship
between
the two
Huffman codes?

If we optimize the
new code,
we also optimize the
old code
(and vice versa).

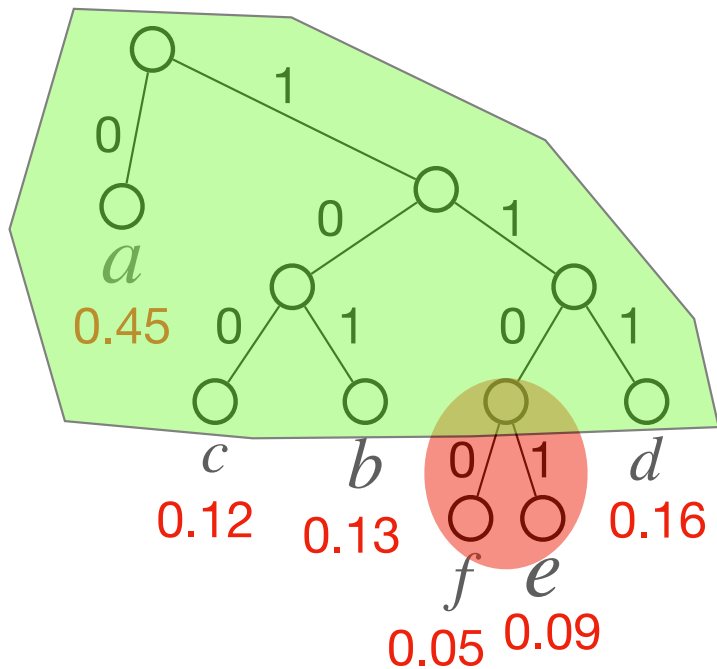
New Code



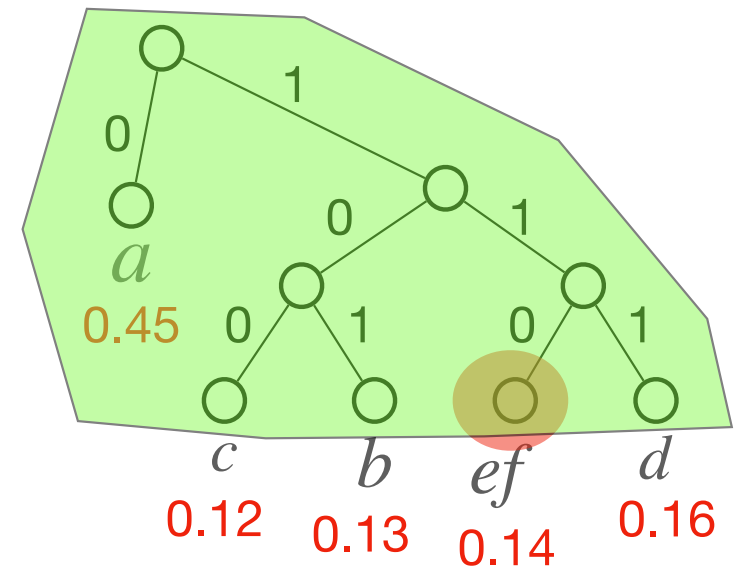
Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>ef</i>
Probability	0.45	0.13	0.12	0.16	0.14

Original Code



New Code



What is
the relationship
between
the two
Huffman codes?

If we optimize the
new code,
we also optimize the
old code
(and vice versa).

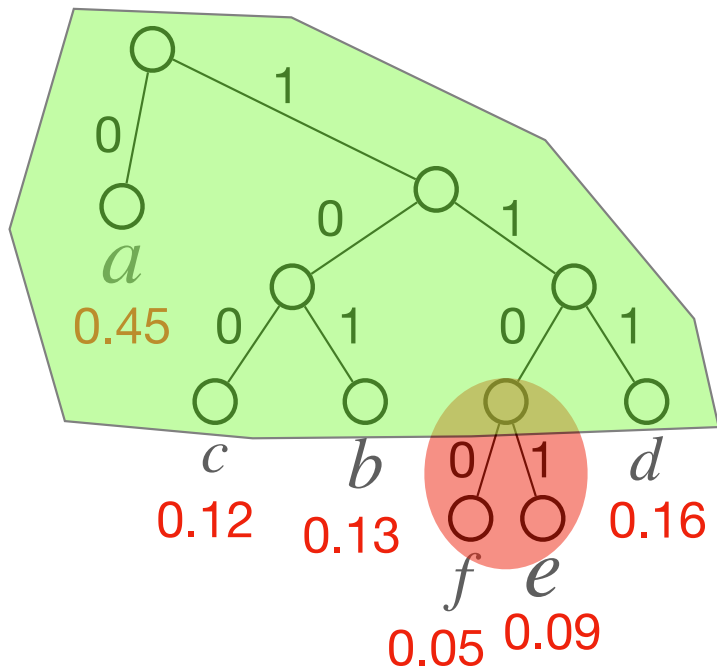
Average codeword length of original code = $1 \times 0.45 + 3 \times 0.13 + 3 \times 0.12 + 3 \times 0.16 + 4 \times 0.09 + 4 \times 0.05$

Average codeword length of new code = $1 \times 0.45 + 3 \times 0.13 + 3 \times 0.12 + 3 \times 0.16 + 3 \times (0.09 + 0.05)$

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>ef</i>
Probability	0.45	0.13	0.12	0.16	0.14

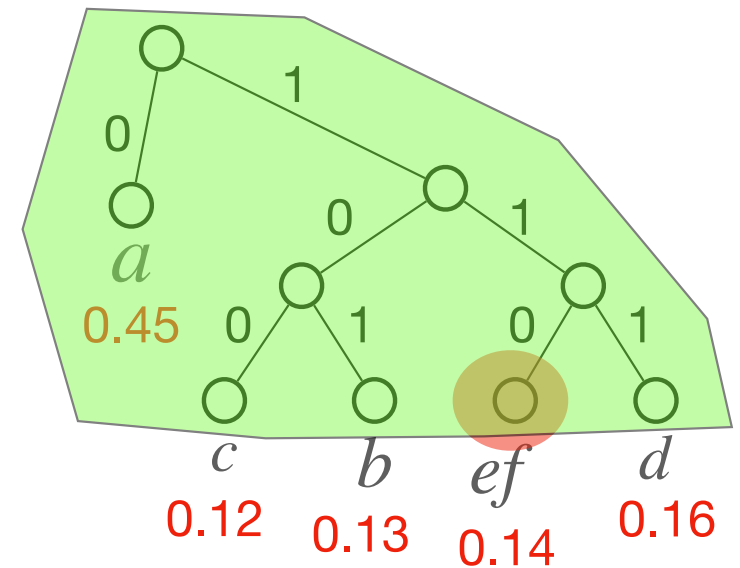
Original Code



What is the relationship between the two Huffman codes?

If we optimize the new code, we also optimize the old code (and vice versa).

New Code



They differ by $f_e + f_f$

Average codeword length of original code = $1 \times 0.45 + 3 \times 0.13 + 3 \times 0.12 + 3 \times 0.16 + 4 \times 0.09 + 4 \times 0.05$

Average codeword length of new code = $1 \times 0.45 + 3 \times 0.13 + 3 \times 0.12 + 3 \times 0.16 + 3 \times (0.09 + 0.05)$

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>ef</i>
Probability	0.45	0.13	0.12	0.16	0.14

Idea of Greedy Algorithm:

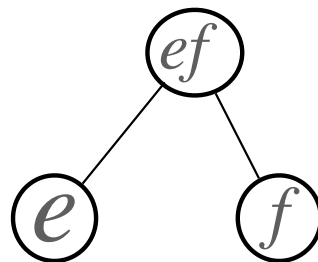
- 1) Make the two symbols of lowest probabilities siblings.
- 2) Combine them into one symbol, and repeat the above process.

Example:

Symbol	a	b	c	d	e	f
Probability	0.45	0.13	0.12	0.16	0.09	0.05

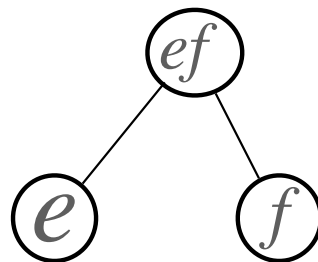
Example:

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>
Probability	0.45	0.13	0.12	0.16	0.09	0.05



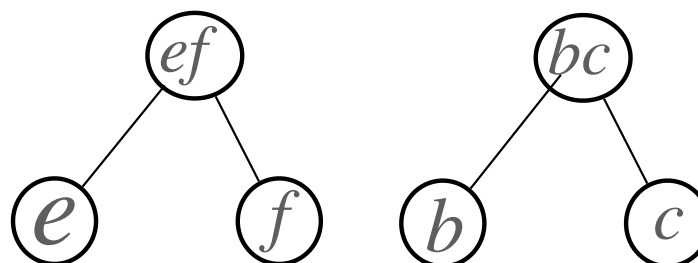
Example:

Symbol	a	b	c	d	ef
Probability	0.45	0.13	0.12	0.16	0.14



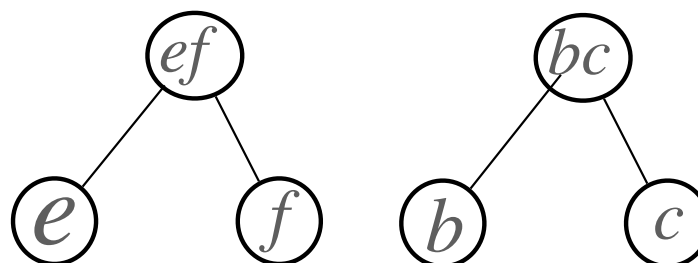
Example:

Symbol	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>ef</i>
Probability	0.45	0.13	0.12	0.16	0.14



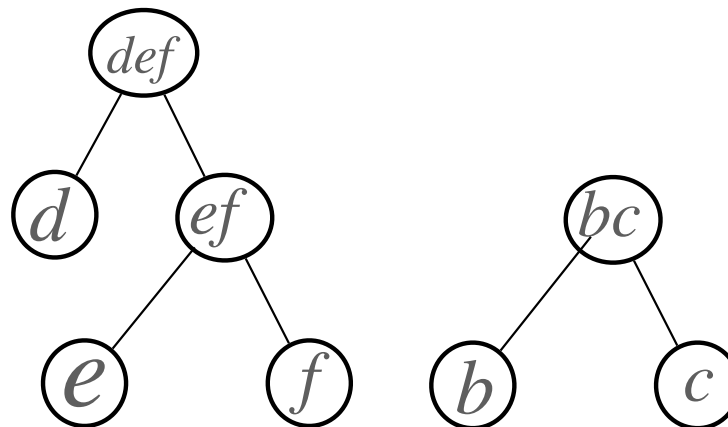
Example:

Symbol	<i>a</i>	<i>bc</i>	<i>d</i>	<i>ef</i>
Probability	0.45	0.25	0.16	0.14



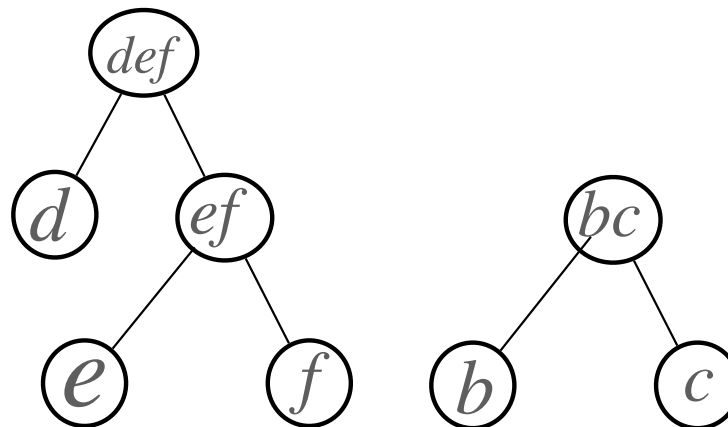
Example:

Symbol	<i>a</i>	<i>bc</i>	<i>d</i>	<i>ef</i>
Probability	0.45	0.25	0.16	0.14



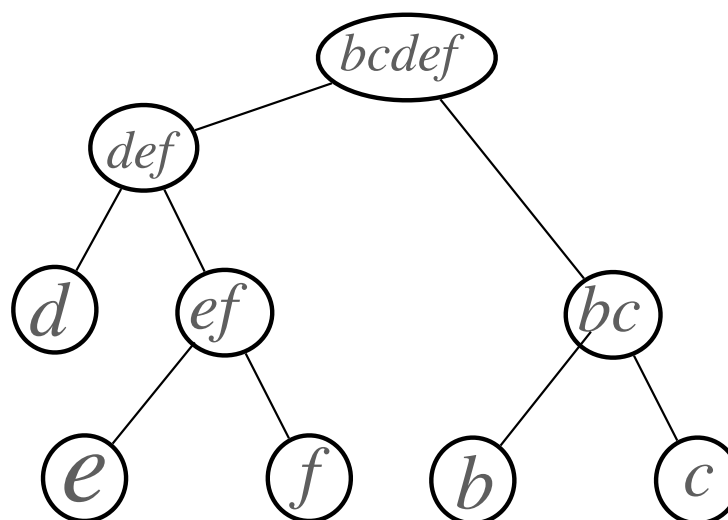
Example:

Symbol	<i>a</i>	<i>bc</i>	<i>def</i>
Probability	0.45	0.25	0.3



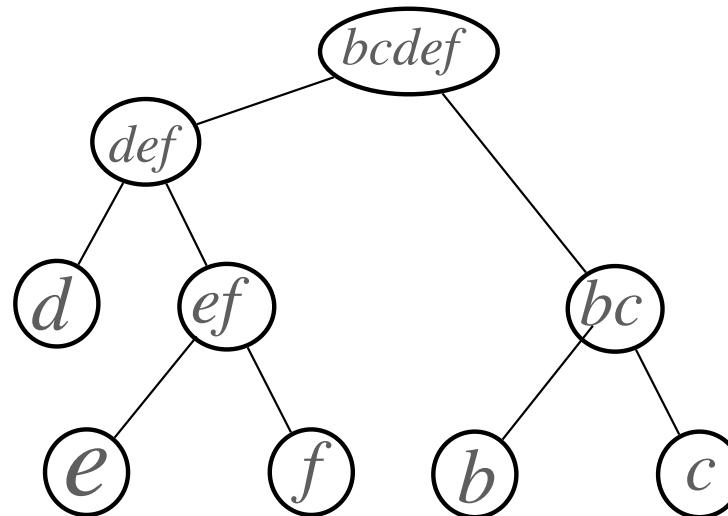
Example:

Symbol	<i>a</i>	<i>bc</i>	<i>def</i>
Probability	0.45	0.25	0.3



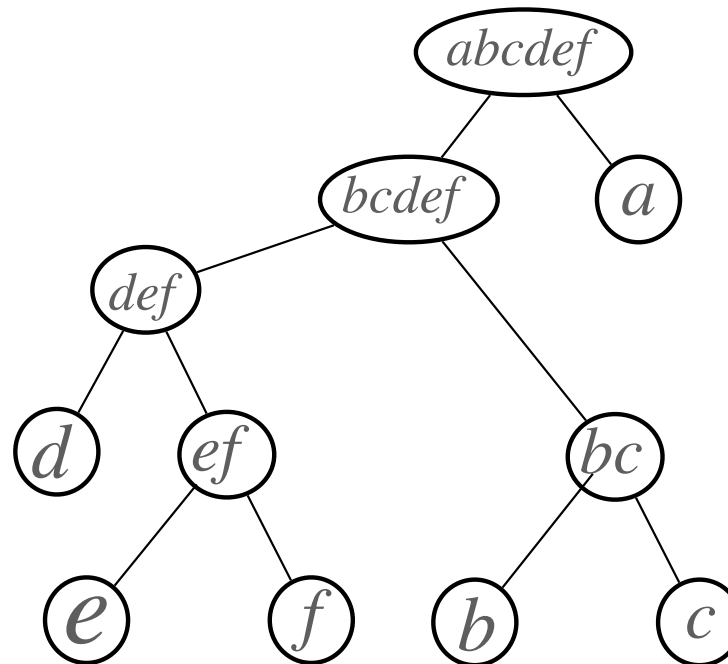
Example:

Symbol	<i>a</i>	<i>bcdef</i>
Probability	0.45	0.55



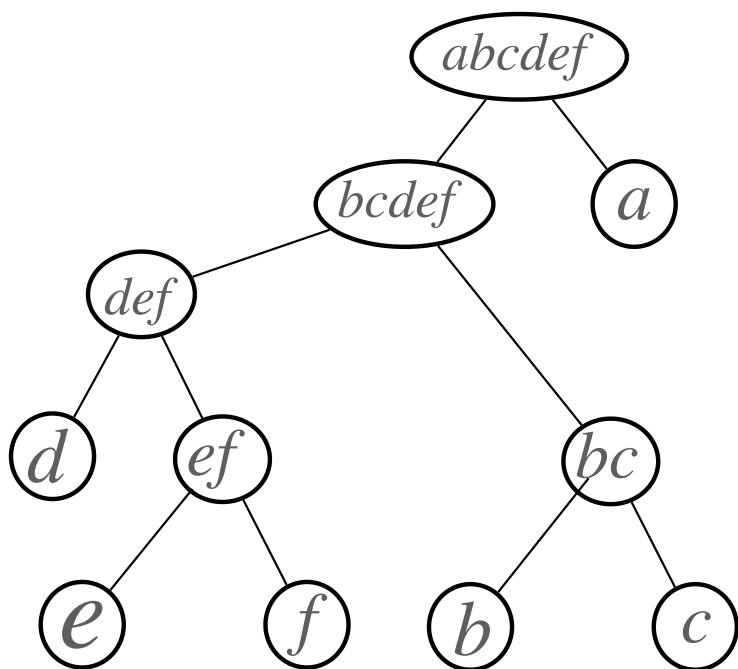
Example:

Symbol	<i>a</i>	<i>bcdef</i>
Probability	0.45	0.55



Example:

Symbol	<i>a</i>	<i>bcdef</i>
Probability	0.45	0.55



Equivalent

