

# High order entropy

Ben Langmead



JOHNS HOPKINS

WHITING SCHOOL  
*of* ENGINEERING

Department of Computer Science



Please sign guestbook ([www.langmead-lab.org/teaching-materials](http://www.langmead-lab.org/teaching-materials)) to tell me briefly how you are using the slides. For original Keynote files, email me ([ben.langmead@gmail.com](mailto:ben.langmead@gmail.com)).

# High-order entropy

Zero order empirical entropy seems insufficient when context matters

Bigram frequency per 40,000 words

<b>th</b>	1.52	<b>en</b>	0.55	<b>ng</b>	0.18
<b>he</b>	1.28	<b>ed</b>	0.53	<b>of</b>	0.16
<b>in</b>	0.94	<b>to</b>	0.52	<b>al</b>	0.09
<b>er</b>	0.94	<b>it</b>	0.50	<b>de</b>	0.09
<b>an</b>	0.82	<b>ou</b>	0.50	<b>se</b>	0.08
<b>re</b>	0.68	<b>ea</b>	0.47	<b>le</b>	0.08
<b>nd</b>	0.63	<b>hi</b>	0.46	<b>sa</b>	0.06
<b>at</b>	0.59	<b>is</b>	0.46	<b>si</b>	0.05
<b>on</b>	0.57	<b>or</b>	0.43	<b>ar</b>	0.04
<b>nt</b>	0.56	<b>ti</b>	0.34	<b>ve</b>	0.04
<b>ha</b>	0.56	<b>as</b>	0.33	<b>ra</b>	0.04
<b>es</b>	0.56	<b>te</b>	0.27	<b>ld</b>	0.02
<b>st</b>	0.55	<b>et</b>	0.19	<b>ur</b>	0.02

# High-order entropy

Can compress better if we consider **context**

Let  $C$  change depending on surrounding symbols

gtgtatcggagcgctctgcgttatcgatcgcgatctggt  
 $C_{tcg}(a)$   $C_{gcg}(a)$   $C_{tcg}(a)$   $C_{tcg}(a)$   $C_{tct}(a)$

For  $k$  symbols of context, we have codes  $C_i \in \{C_{\Sigma^k}\}$

# High-order entropy

Could consider context to the **right** →

gtgt**at****c**gg**a**gc**g**ctctg**c**gtt**at****c**g**a**t**c**g**c**g**a**t**c**tgggt  
 $C_{tcg}(a)$   $C_{gcg}(a)$   $C_{tcg}(a)$   $C_{tcg}(a)$   $C_{tct}(a)$

Or context to the **left** ←

gtgt**a**t**c**gg**a**gc**g**ctctg**c**gtt**a**t**c**g**a**t**c**g**c**g**a**t**c**tgggt  
 $C_{tct}(a)$   $C_{cgg}(a)$   $C_{ggt}(a)$   $C_{tcg}(a)$   $C_{gcg}(a)$

# High-order entropy

How should we build each code  $\{C_{\Sigma^k}\}$ ?

Same as before, but with frequencies  
*conditioned on context*

E.g.  $C_{gca}$  is built considering the number of  
times each symbol occurs just after gca

# High-order entropy

abracadabraabracadabra

Let  $S_a$  be the substring we get by concatenating characters just after the a's

$$S_a = bcd b a b c d b$$

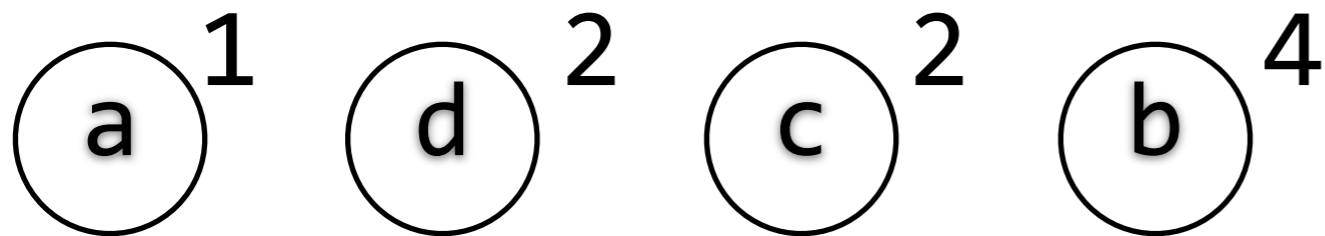
Build code  $C_a$  using frequencies in  $S_a$ :

{a : 1, b : 4, c : 2, d : 2, r : 0}

(r won't get a code)

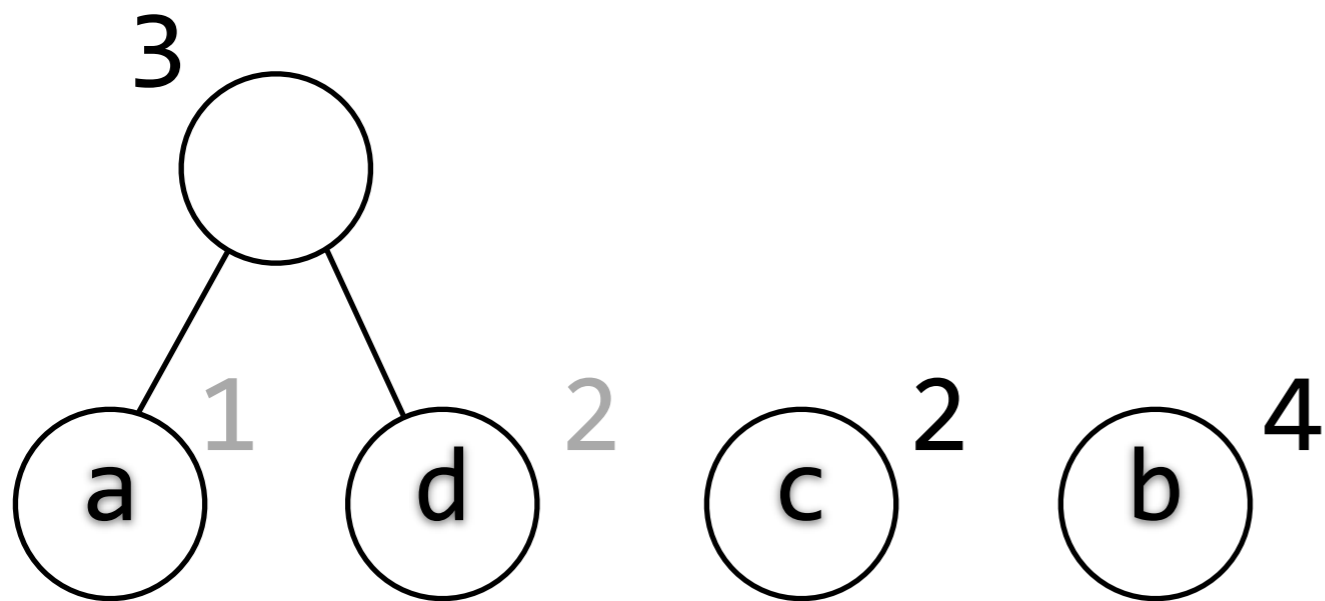
# High-order entropy

{a : 1, b : 4, c : 2, d : 2, r : 0}



# High-order entropy

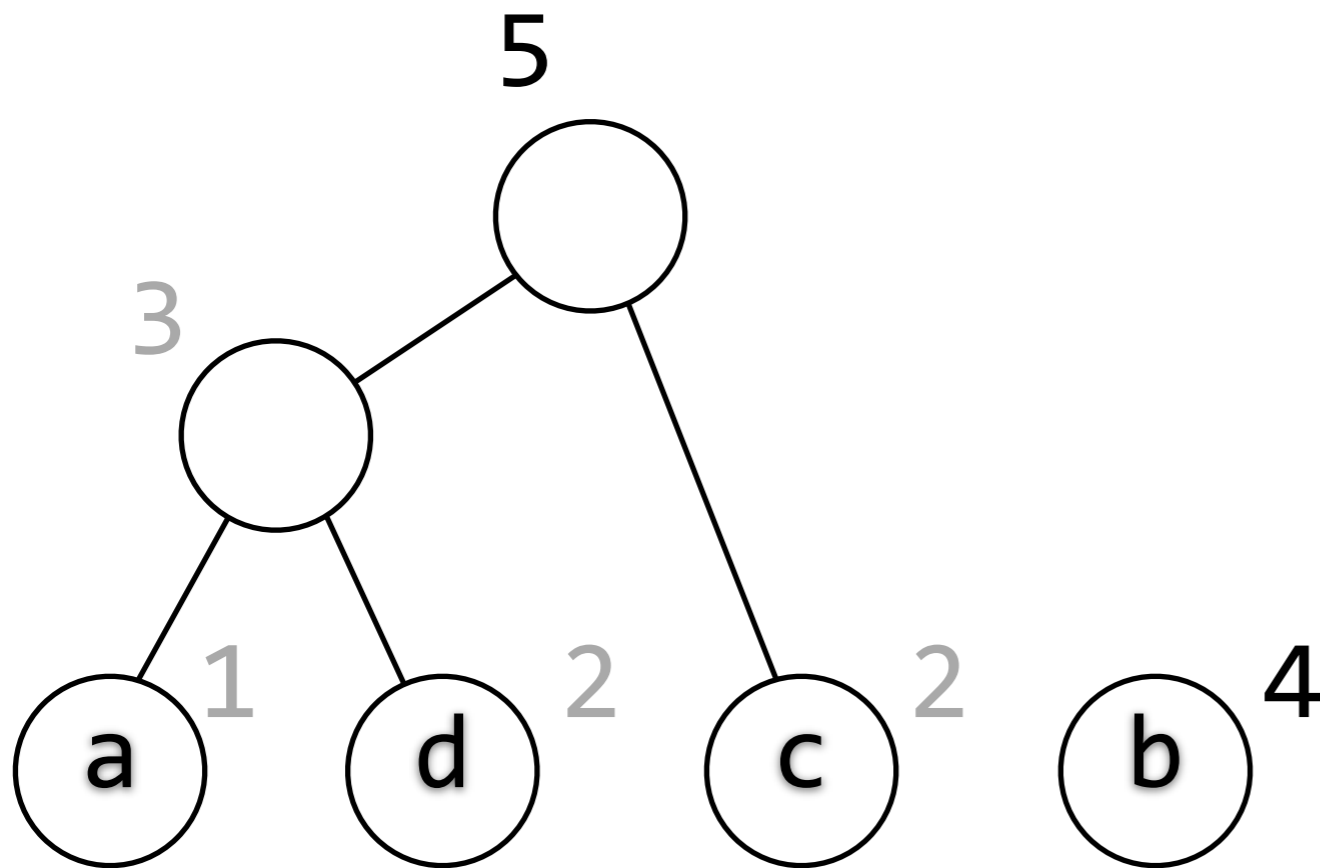
$\{a : 1, b : 4, c : 2, d : 2, r : 0\}$





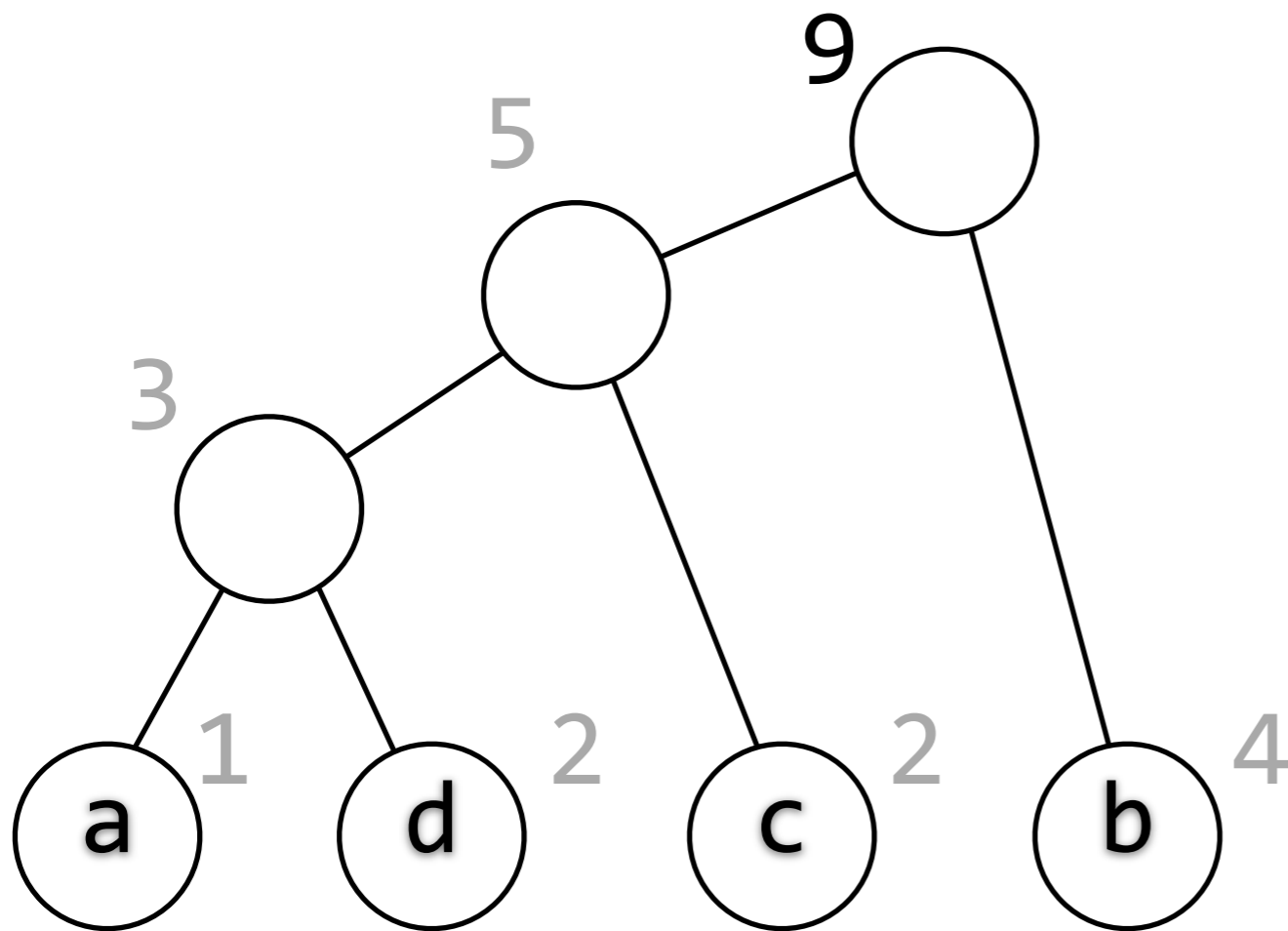
# High-order entropy

$\{a : 1, b : 4, c : 2, d : 2, r : 0\}$



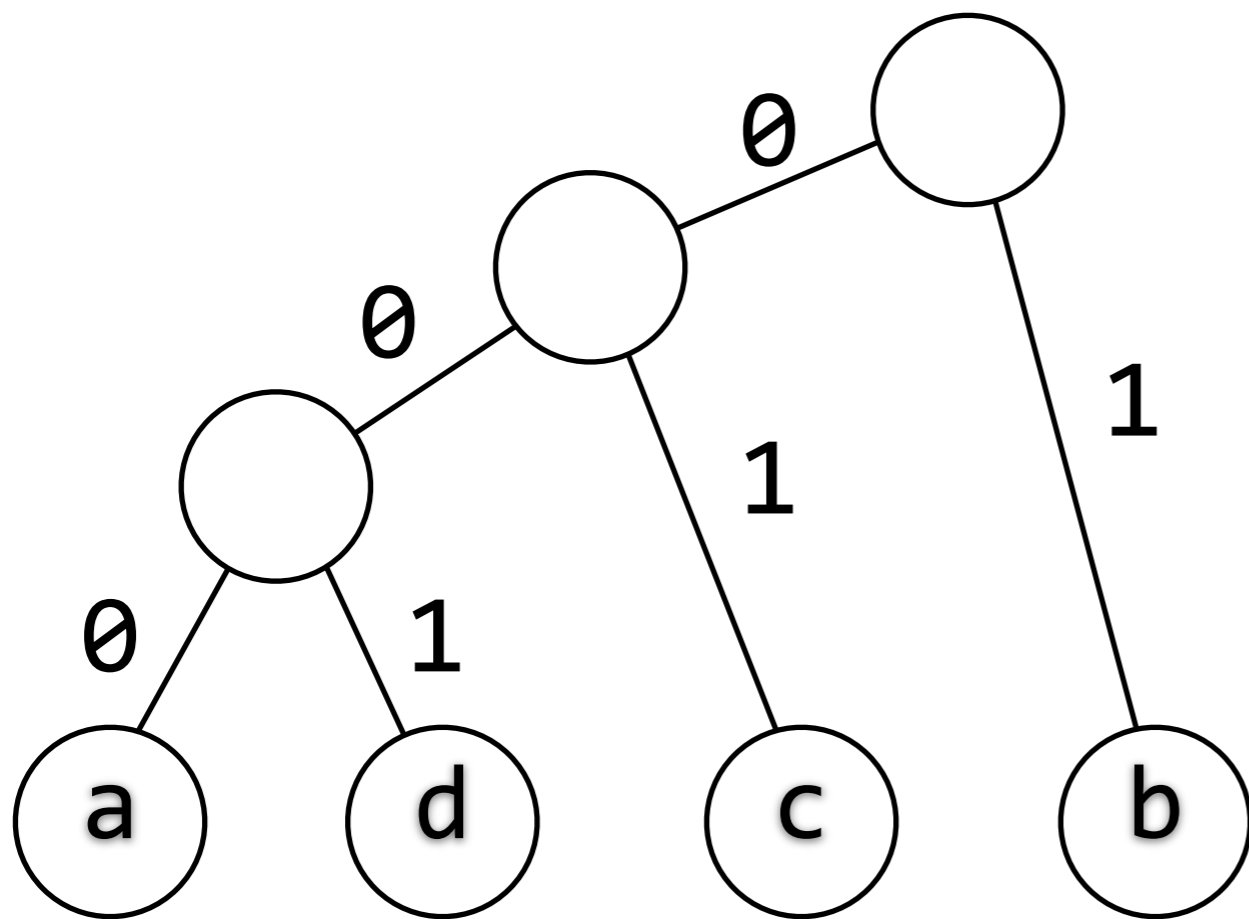
# High-order entropy

{a : 1, b : 4, c : 2, d : 2, r : 0}



# High-order entropy

{a : 1, b : 4, c : 2, d : 2, r : 0}



$$C_a(a) = 000$$

$$C_a(d) = 001$$

$$C_a(c) = 01$$

$$C_a(b) = 1$$

# High-order entropy

abracadabraabracadabra

$$S_a = \text{bcdbabcdb}$$

$$\{a : 1, b : 4, c : 2, d : 2\}$$

$$C_a(a) = 000$$

$$C_a(d) = 001$$

$$C_a(c) = 01$$

$$C_a(b) = 1$$

$$S_b = \text{rrrrr} \quad \{r : 4\}$$

(no code)

$$S_c = \text{aa} \quad \{a : 2\}$$

(no code)

$$S_d = \text{aa} \quad \{a : 2\}$$

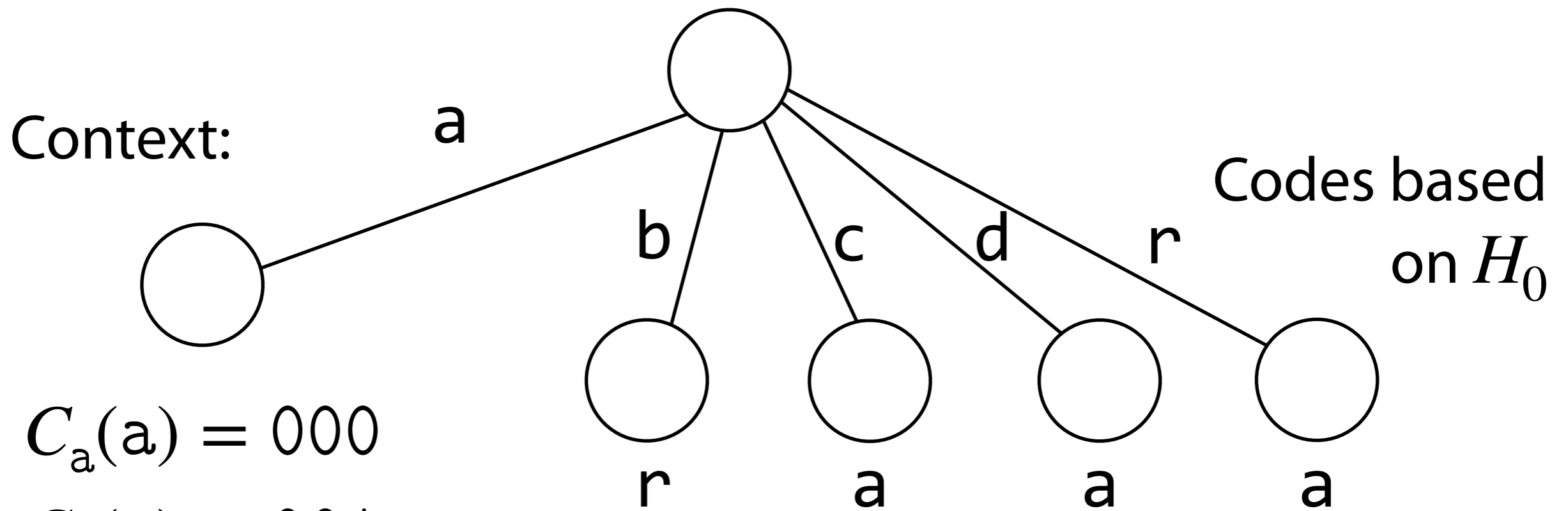
(no code)

$$S_r = \text{aaaa} \quad \{a : 4\}$$

(no code)

# High-order entropy

abracadabraabracadabra



$$C_a(a) = 000$$

$$C_a(b) = 001$$

$$C_a(c) = 01$$

$$C_a(d) = 1$$

(no codes required)

# High-order entropy

mississippi

$S_i = \text{sspmsp}$

$\{s : 4, p : 2, m : 1\}$

$C_i(p) = 00$

$C_i(m) = 01$

$C_i(s) = 1$

$S_m = \text{ii} \quad \{i : 2\}$

(no code)

$S_p = \text{pipi}$

$C_p(p) = 0$

$\{p : 2, i : 2\}$

$C_p(i) = 1$

$S_s = \text{sisisisi}$

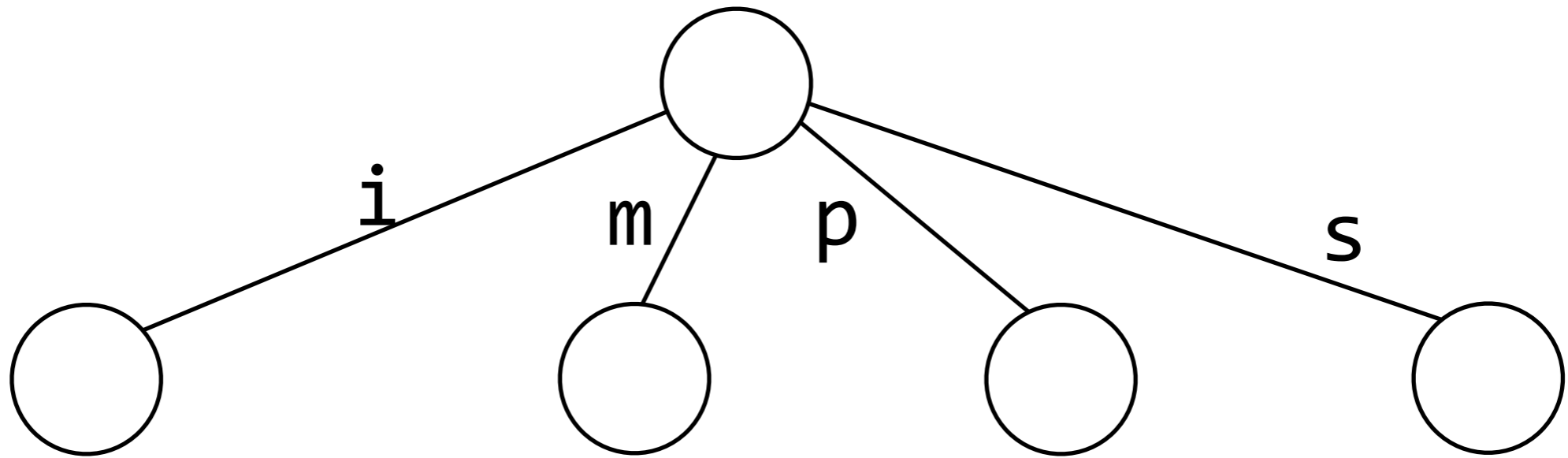
$\{s : 4, i : 4\}$

$C_s(s) = 0$

$C_s(i) = 1$

# High-order entropy

mississippiissippi



$$C_i(p) = 00$$

$$C_i(m) = 01$$

$$C_i(s) = 1$$

i

(no code)

$$C_p(p) = 0$$

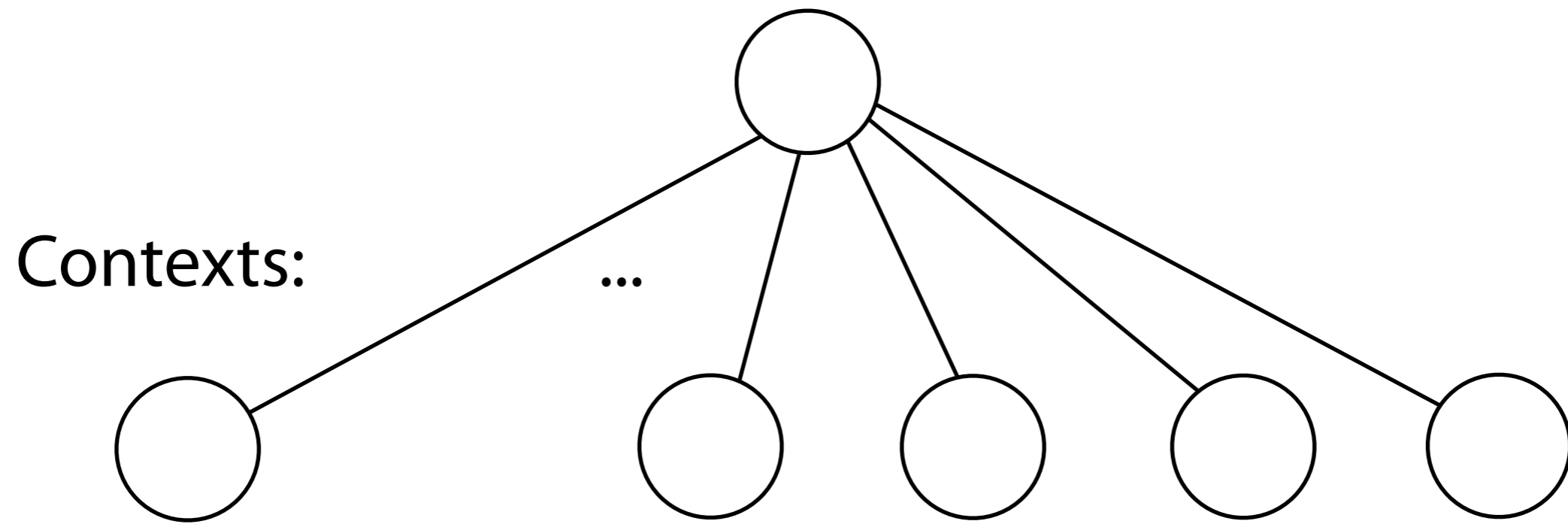
$$C_p(i) = 1$$

$$C_s(s) = 0$$

$$C_s(i) = 1$$

# High-order entropy

Def'n of **high-order empirical entropy**  $H_k$  is similarly hierarchical



Codes achieving near- $H_0$  given context



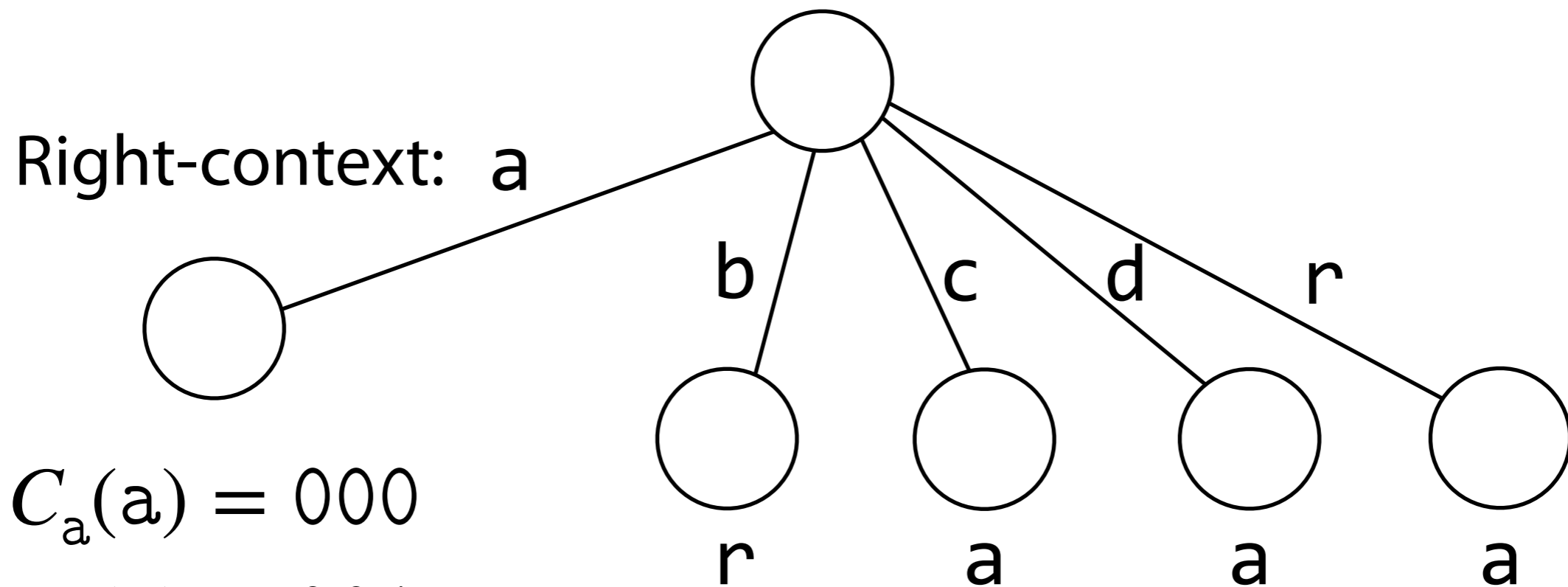
# High-order entropy

$H_k$  of a length- $n$  string  $S$  is a weighted sum **over all contexts** of the **zero order empirical entropy** of symbols having that context

$$H_k(S) = \sum_{t \in \Sigma^k} \frac{|S_t|}{n} \cdot H_0(S_t) \quad \text{for } k > 0$$

$S$  is the entire string,  $S_t$  is the concatenation of symbols having context  $t$

# High-order entropy



$$C_a(\mathbf{a}) = 000$$

$$C_a(\mathbf{d}) = 001$$

$$C_a(\mathbf{c}) = 01$$

$$C_a(\mathbf{b}) = 1$$

Schemes like this can compress  
to  $\leq n(H_k(S) + 1)$  bits

With added overhead of  
switching between many codes

# High-order entropy

Collection	H0	H1	H2
CODE SOURCES	5.537 (69.21%)	4.038 (50.48%)	3.012 (37.65%)
MIDI	5.633 (70.41%)	4.734 (59.18%)	4.139 (51.74%)
PROTEINS	4.195 (52.44%)	4.173 (52.16%)	4.146 (51.82%)
DNA	1.982 (24.78%)	1.935 (24.19%)	1.925 (24.06%)
ENGLISH	4.529 (56.61%)	3.606 (45.08%)	2.922 (36.53%)
XML	5.230 (65.37%)	3.294 (41.17%)	2.007 (25.09%)

Empirical entropies for 6 texts. Values are bits-per-symbol, percentages are ratios compared to ASCII.

Table from: Prezza, Nicola. Compressed Computation for Text Indexing. Diss. PhD thesis, University of Udine, 2016.

# High-order entropy

$H_k$  encoding reaches into the string, extracting "structure" needed to compress well

$k$  balances compression with overhead

*Grouping* principle at play

$H_0$ -based methods are simpler, faster, require less memory, but can't find as much structure as  $H_k$

...or...can they?

*Order* to the rescue