

High order entropy

Ben Langmead



JOHNS HOPKINS

WHITING SCHOOL
of ENGINEERING

Department of Computer Science



Please sign guestbook (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).

High-order entropy

Zero order empirical entropy seems insufficient when context matters

Bigram frequency per 40,000 words

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

High-order entropy

Can compress better if we consider **context**

Let C change depending on surrounding symbols

gtgt**a**tcgg**a**gcgctctgcgtt**a**tcg**a**tcg**a**tcg**a**ctgggt
 $C_{tcg}(a)$ $C_{gcg}(a)$ $C_{tcg}(a)$ $C_{tcg}(a)$ $C_{tcg}(a)$ $C_{tct}(a)$

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**cgga**gc**gctctg**cg**tt**at**cg**at**cg**cg****at**ctggta
 $C_{tcg}(a)$ $C_{gcg}(a)$ $C_{tcg}(a)$ $C_{tcg}(a)$ $C_{tct}(a)$

Or context to the **left** ←

gtgt**at**cgga**gc**gctctg**cg**tt**at**cg**at**cg**cg****at**ctggta
 $C_{tct}(a)$ $C_{cgg}(a)$ $C_{ggt}(a)$ $C_{tcg}(a)$ $C_{gcg}(a)$

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 string we get by concatenating characters just after the a's

$$S_a =$$

High-order entropy

abracadabraabracadabra

Let S_a be the string 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 :

High-order entropy

abracadabraabracadabra

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

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

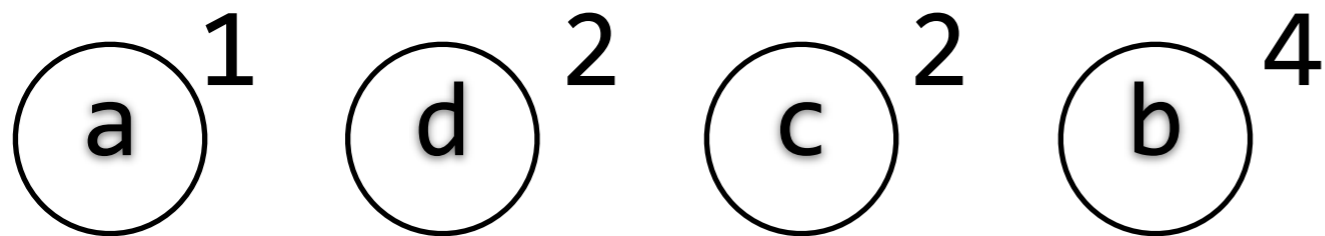
Build code C_a using frequencies in S_a :

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

(r won't get a code)

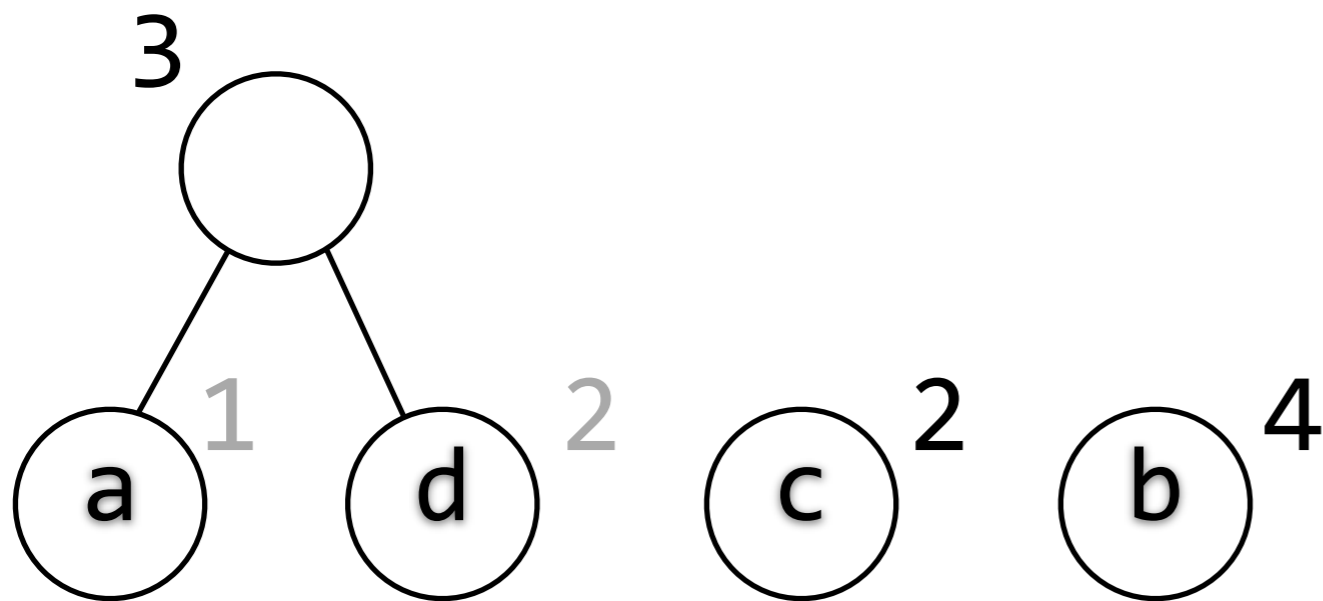
High-order entropy

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



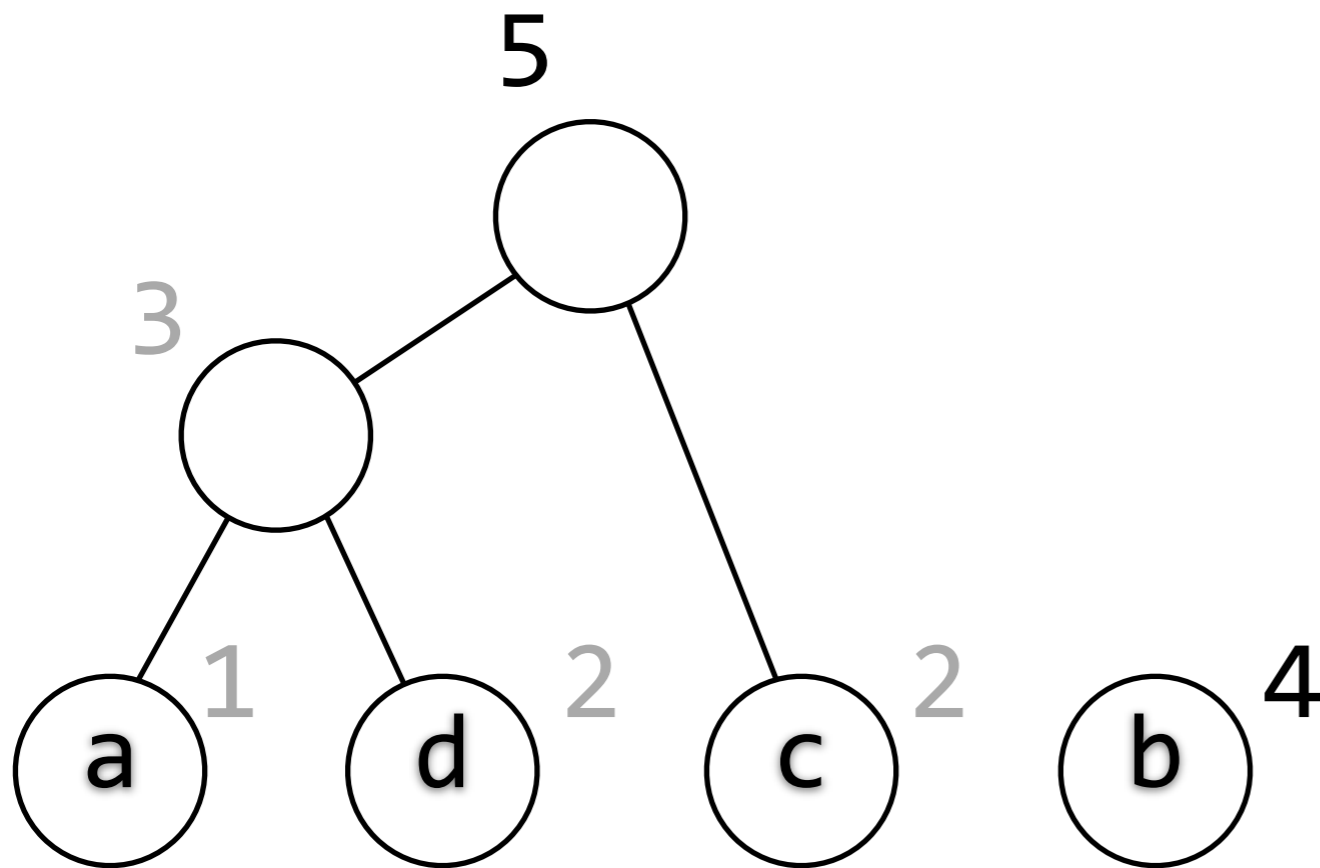
High-order entropy

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



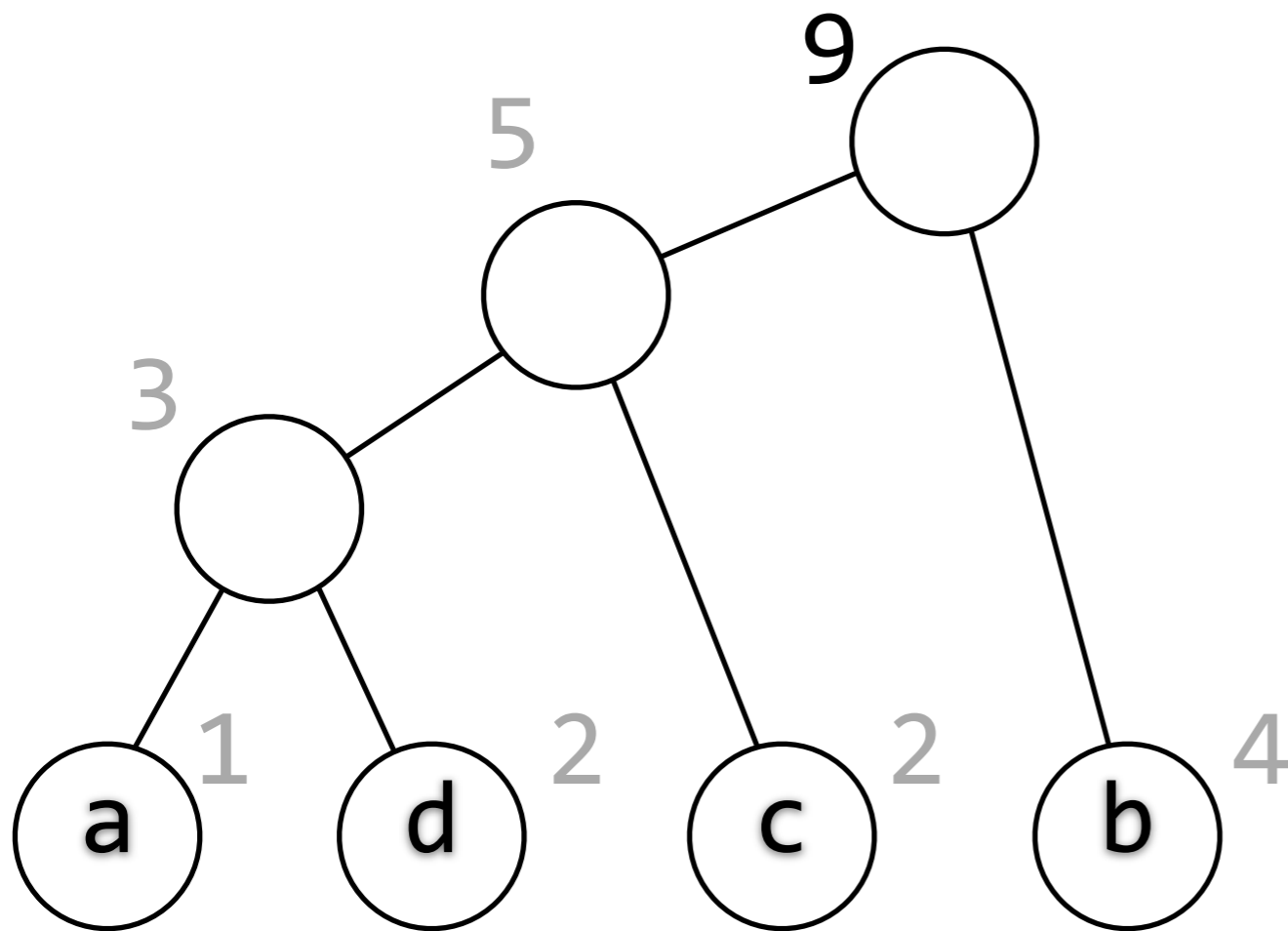
High-order entropy

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



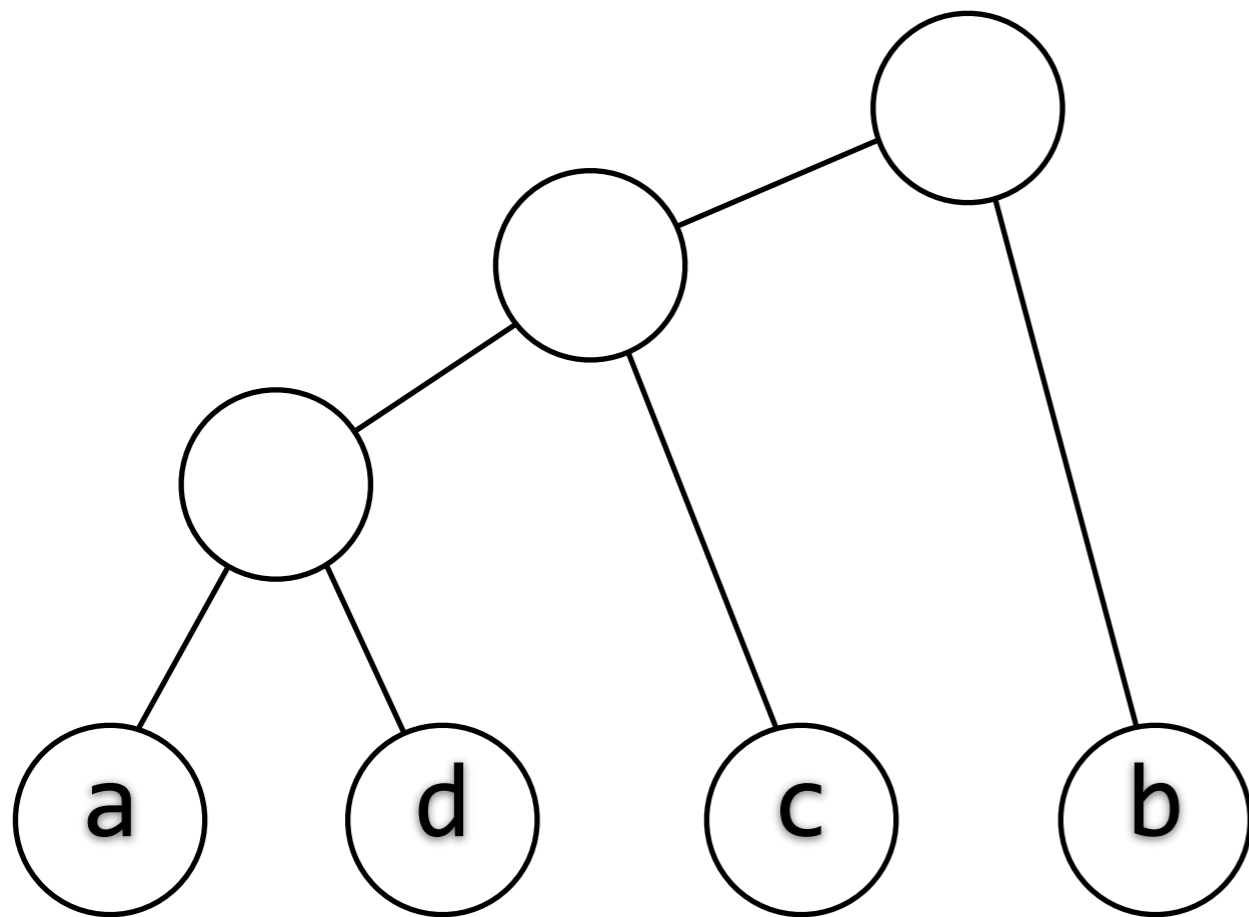
High-order entropy

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



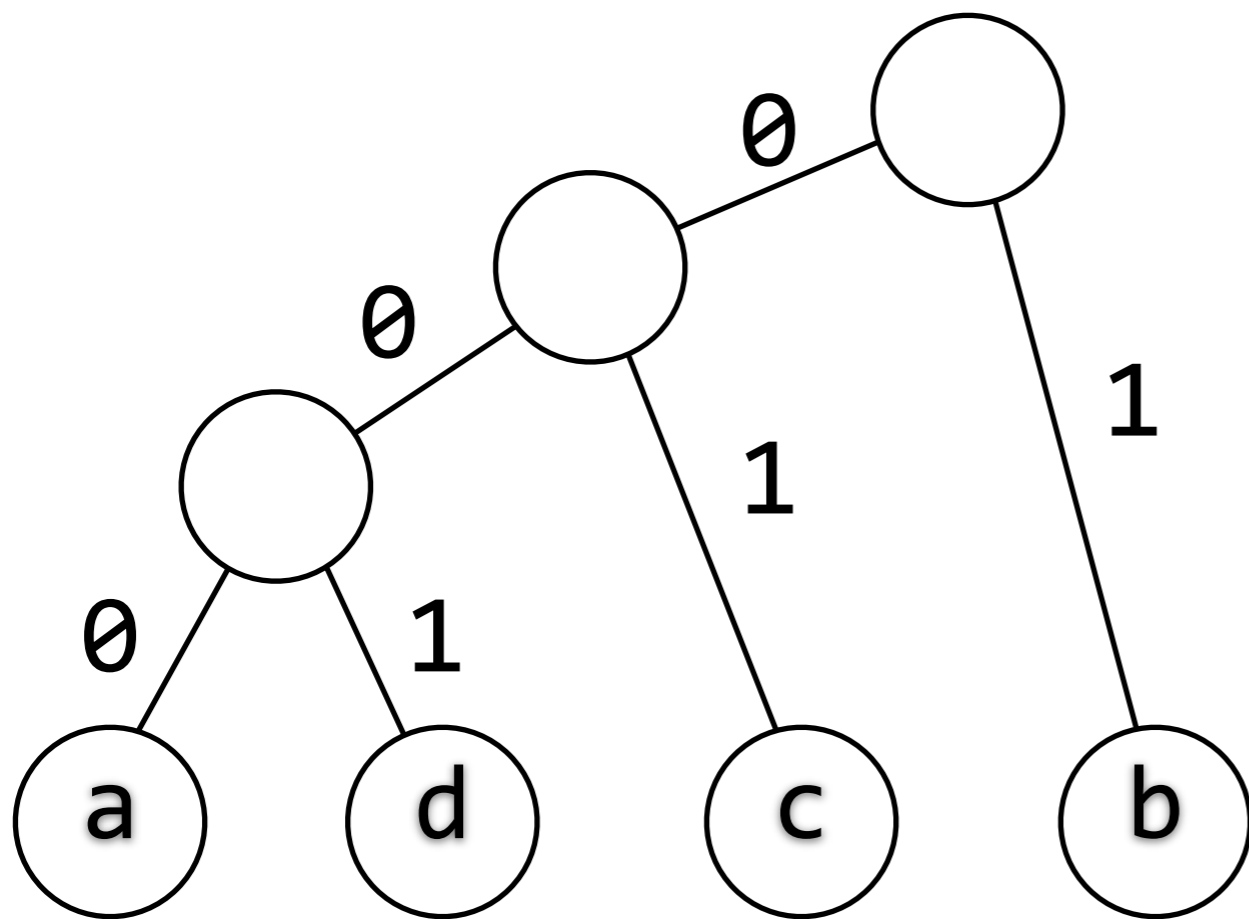
High-order entropy

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



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

mississippi

$S_i = \text{sspmspp}$

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

$C_i(p) = 00$

$C_i(m) = 01$

$C_i(s) = 1$

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

(no code)

$S_p = \text{pipi}$

$C_p(p) = 0$

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

$C_p(i) = 1$

$S_s = \text{sisisisi}$

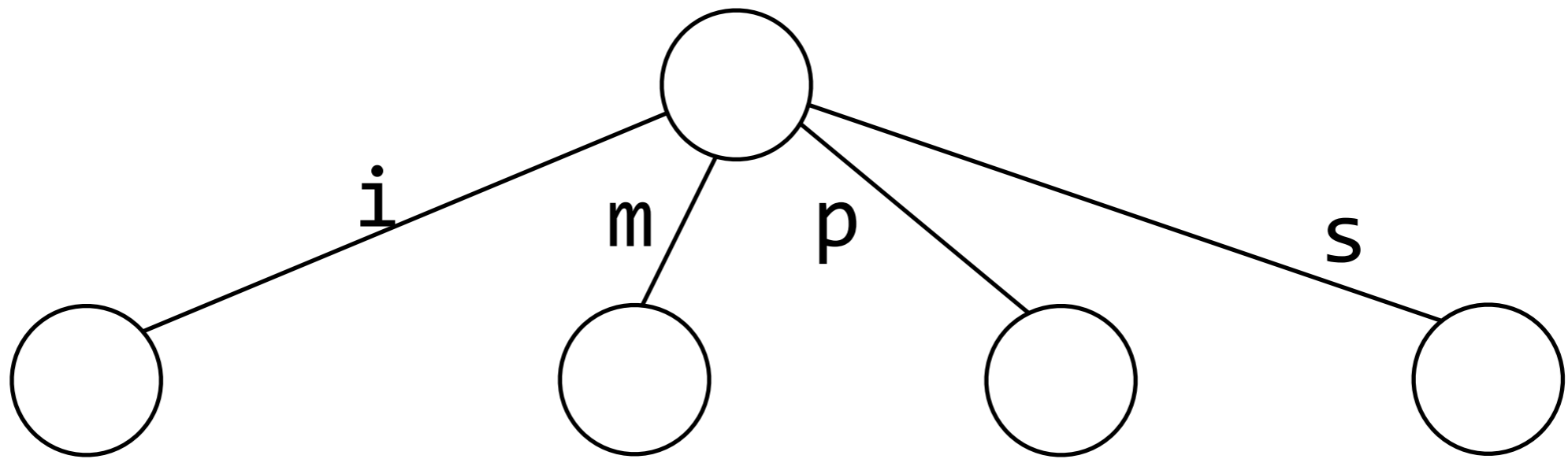
$\{s : 4, i : 4\}$

$C_s(s) = 0$

$C_s(i) = 1$

High-order entropy

mississippiissippi



$$C_i(\mathbf{p}) = 00$$

$$C_i(\mathbf{m}) = 01$$

$$C_i(\mathbf{s}) = 1$$

i

(no code)

$$C_p(\mathbf{p}) = 0$$

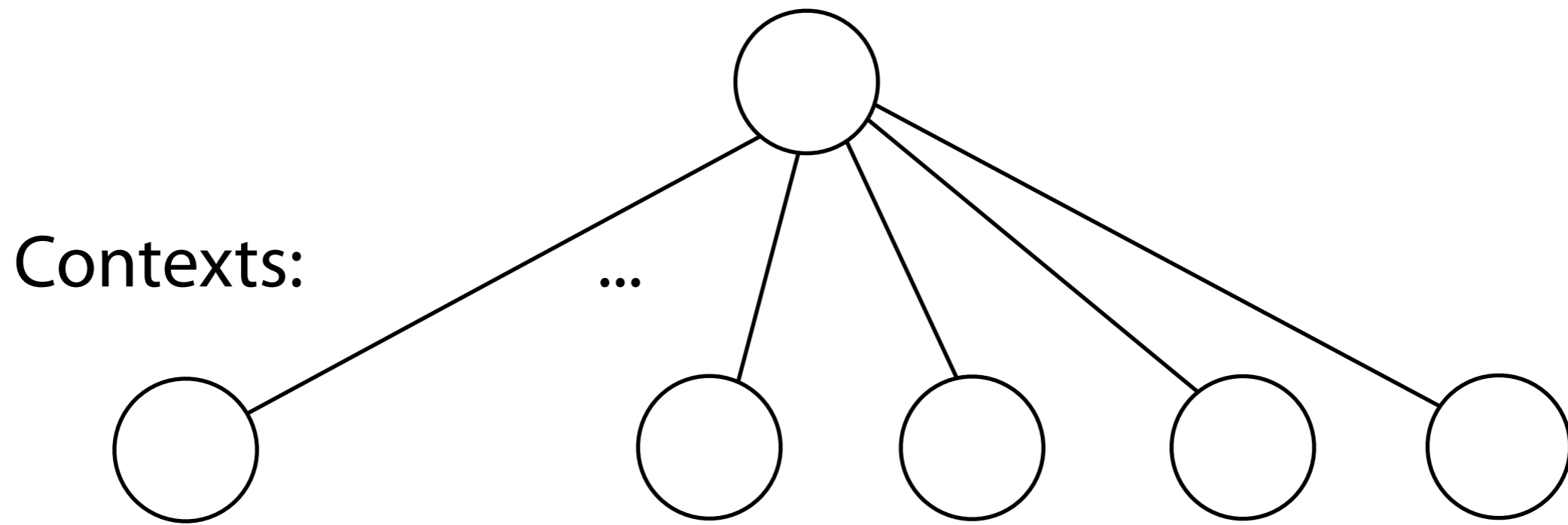
$$C_p(\mathbf{i}) = 1$$

$$C_s(\mathbf{s}) = 0$$

$$C_s(\mathbf{i}) = 1$$

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

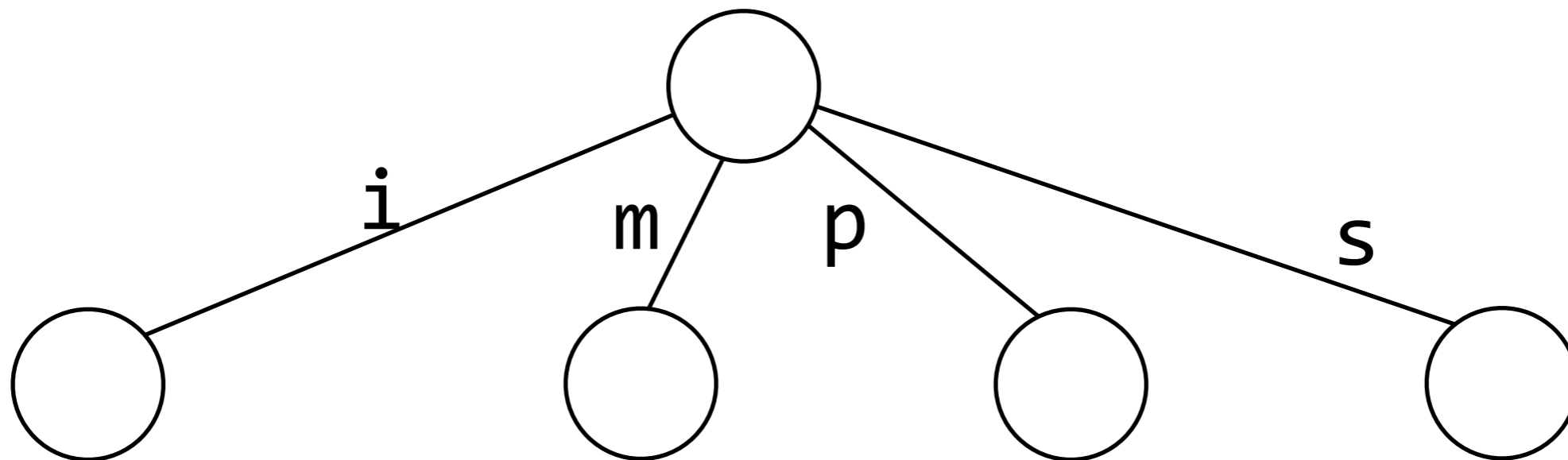
High-order entropy

H_k of 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_t is the concatenation of symbols having context t

High-order entropy



$$C_i(\mathbf{p}) = 00$$

$$C_i(\mathbf{m}) = 01$$

$$C_i(\mathbf{s}) = 1$$

\mathbf{i}
(no code)

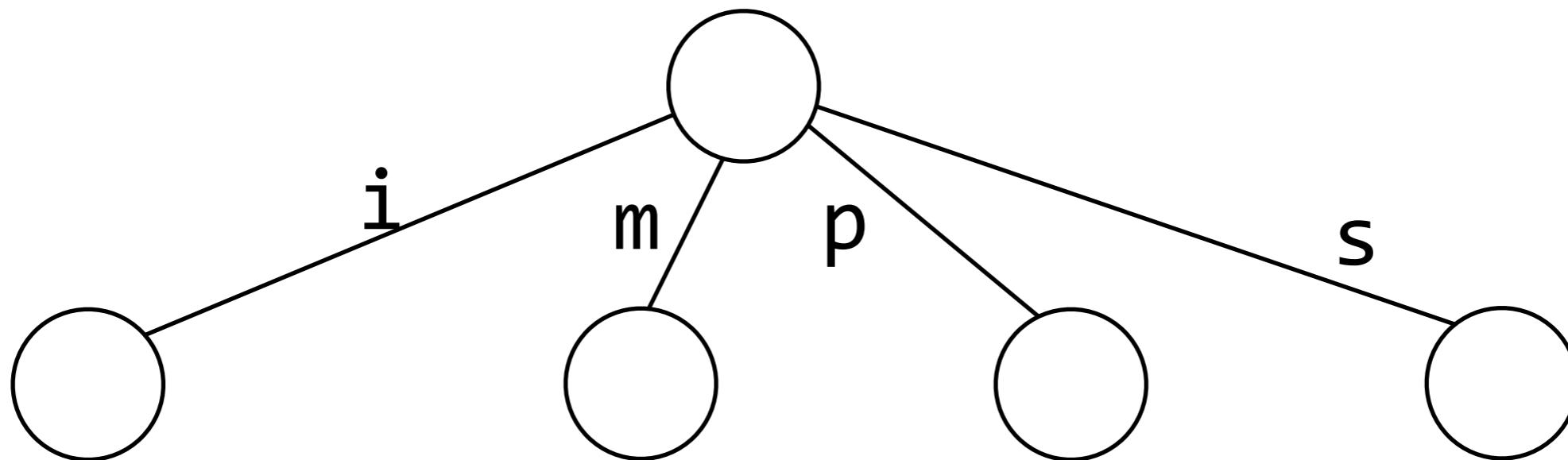
$$C_p(\mathbf{p}) = 0$$

$$C_p(\mathbf{i}) = 1$$

$$C_s(\mathbf{s}) = 0$$

$$C_s(\mathbf{i}) = 1$$

High-order entropy



$$C_i(\mathbf{p}) = 00$$

$$C_i(\mathbf{m}) = 01$$

$$C_i(\mathbf{s}) = 1$$

\mathbf{i}

(no code)

$$C_p(\mathbf{p}) = 0$$

$$C_p(\mathbf{i}) = 1$$

$$C_s(\mathbf{s}) = 0$$

$$C_s(\mathbf{i}) = 1$$

Can compress to
 $\leq n(H_k(S) + 1)$ bits

Switching between many
codes adds overhead

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