

Tries

("try"-s)

CMSC 420

A brief history of character sets

- 1963: American Standard Code for Information Interchange (**ASCII**)
 - Used numbers from 32 through 127 to encode *unaccented* English characters and some common symbols, like >, <, _, etc
 - The numbers from 0 through 31 were used for **control characters** that are not all particularly useful nowadays, like BELL, printer feed, etc.

A brief history of character sets

- 1963: American Standard Code for Information Interchange (**ASCII**)
 - Used numbers from 32 through 127 to encode *unaccented* English characters and some common symbols, like >, <, _, etc
 - The numbers from 0 through 31 were used for **control characters** that are not all particularly useful nowadays, like BELL, printer feed, etc.
- Need only $\log_2 128 = 7$ bits to encode those, so with 1 extra bit, you can do whatever you want!
- Example: the IBM-PC, which was using an Intel 8088 CPU used a standard known as OEM...

OEM standard

F * ? / ? O - o d g f » T ■ n
E □ A . > N ^ n ^ a R < F ^ H E ■
D + - = M] e x i y i u = - s z
C ♀ , < L \ l i t E x u H e n
B δ ← + ; K [k C i c x H δ
A → * : J Z j N e ü R || - G .
↓ > 9 I Y i u e ö R H G B .
↑ < 8 H X h x e ü c R H G B .
‡ , 7 G W g w c u e || - x
6 ♠ = & 6 F U f u a u d H F p u .
5 ♦ ♣ x 5 E U e u à o ñ F J
4 ♦ ♣ \$ 4 D T d t ä ö H - E
3 ♥ :: # 3 C S c s a ö u - H E
2 ♣ ♦ " 2 B R b r é F ö T P R
1 ♣ ♦ ! 1 A Q a q ü a i T P
0 ♣ ♦ 0 @ P , p c E á L u x
0 1 2 3 4 5 6 7 8 9 A B C D E F

OEM standard

0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	®	♥	♦	♠	♣	—	↑	↓	→	←	↶	↷	▀	*	
1	►	◀	↕	॥	₪	₭	₩	₱	₭	₪	₭	₵	₭	₭	₭
2	:	”	”	#	\$	₹	&	,	<	>	*	+	,	-	/
3	Ø	1	2	3	4	5	6	7	8	9	:	;	<	=	>?
4	Œ	À	Ã	Ɓ	Ҫ	Ɗ	Ӗ	Ӯ	ӻ	ӱ	Ӱ	Ӳ	Ӵ	ӵ	Ӷ
5	܂	܃	܄	܅	܆	܇	܈	܉	܊	܋	܌	܍	܎	܏	ܐ
6	܂	܃	܄	܅	܆	܇	܈	܉	܊	܋	܌	܍	܎	܏	ܐ
7	܂	܃	܄	܅	܆	܇	܈	܉	܊	܋	܌	܍	܎	܏	ܐ
8	܂	܃	܄	܅	܆	܇	܈	܉	܊	܋	܌	܍	܎	܏	ܐ
9	܂	܃	܄	܅	܆	܇	܈	܉	܊	܋	܌	܍	܎	܏	ܐ
A	܂	܃	܄	܅	܆	܇	܈	܉	܊	܋	܌	܍	܎	܏	ܐ
B	܂	܃	܄	܅	܆	܇	܈	܉	܊	܋	܌	܍	܎	܏	ܐ
C	܂	܃	܄	܅	܆	܇	܈	܉	܊	܋	܌	܍	܎	܏	ܐ
D	܂	܃	܄	܅	܆	܇	܈	܉	܊	܋	܌	܍	܎	܏	ܐ
E	܂	܃	܄	܅	܆	܇	܈	܉	܊	܋	܌	܍	܎	܏	ܐ
F	܂	܃	܄	܅	܆	܇	܈	܉	܊	܋	܌	܍	܎	܏	ܐ

Last ASCII character,
since $127_{(10)} = 0x7E$

OEM standard

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	☺	☻	♥	♦	♣	♠					δ	♀	♫	*		
1	►	◀	↕	॥	፩	፪	=	↑	↓	→	←	↶	↷	▲	▼	
2	:	"	#	\$	٪	&	,	<	>	*	+	,	-	.	/	
3	፦	፧	፪	፫	፬	፭	፮	፯	፱	፻	፼	፽	፾	፷	፸	
4	፻	፳	፴	፵	፶	፷	፸	፹	፺	፻	፼	፽	፿	፻	፻	
5	፻	፳	፴	፵	፶	፷	፸	፹	፺	፻	፼	፽	፿	፻	፻	
6	፻	፳	፴	፵	፶	፷	፸	፹	፺	፻	፼	፽	፿	፻	፻	
7	፻	፳	፴	፵	፶	፷	፸	፹	፺	፻	፼	፽	፿	፻	፻	
8	፻	፳	፴	፵	፶	፷	፸	፹	፺	፻	፼	፽	፿	፻	፻	
9	፻	፳	፴	፵	፶	፷	፸	፹	፺	፻	፼	፽	፿	፻	፻	
A	፻	፳	፴	፵	፶	፷	፸	፹	፺	፻	፼	፽	፿	፻	፻	
B	፻	፳	፴	፵	፶	፷	፸	፹	፺	፻	፼	፽	፿	፻	፻	
C	፻	፳	፴	፵	፶	፷	፸	፹	፺	፻	፼	፽	፿	፻	፻	
D	፻	፳	፴	፵	፶	፷	፸	፹	፺	፻	፼	፽	፿	፻	፻	
E	፻	፳	፴	፵	፶	፷	፸	፹	፺	፻	፼	፽	፿	፻	፻	
F	፻	፳	፴	፵	፶	፷	፸	፹	፺	፻	፼	፽	፿	፻	፻	

ASCII

I'm the IBM-PC and I
do what I want with
my top bit

OEM standard

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	☺	☻	♥	♦	♣	♠	⌚	⌚	⌚	⌚	⌚	⌚	⌚	⌚	⌚	
1	◀	▶	↔	☰	✖	✖	✖	✖	✖	✖	✖	✖	✖	✖	✖	
2	?	"	#	\$	%	&	,	<	>	*	+	,	-	/		
3	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>?	
4	Ѐ	Ѐ	Ѐ	Ѐ	Ѐ	Ѐ	Ѐ	Ѐ	Ѐ	Ѐ	Ӣ	Ӣ	Ӣ	Ӣ	Ӣ	
5	܂	܂	܂	܂	܂	܂	܂	܂	܂	܂	܂	܂	܂	܂	܂	
6	܃	܃	܃	܃	܃	܃	܃	܃	܃	܃	܃	܃	܃	܃	܃	
7	܄	܄	܄	܄	܄	܄	܄	܄	܄	܄	܄	܄	܄	܄	܄	
8	܅	܅	܅	܅	܅	܅	܅	܅	܅	܅	܅	܅	܅	܅	܅	
9	܆	܆	܆	܆	܆	܆	܆	܆	܆	܆	܆	܆	܆	܆	܆	
A	܇	܇	܇	܇	܇	܇	܇	܇	܇	܇	܇	܇	܇	܇	܇	
B	܈	܈	܈	܈	܈	܈	܈	܈	܈	܈	܈	܈	܈	܈	܈	
C	܉	܉	܉	܉	܉	܉	܉	܉	܉	܉	܉	܉	܉	܉	܉	
D	܊	܊	܊	܊	܊	܊	܊	܊	܊	܊	܊	܊	܊	܊	܊	
E	܋	܋	܋	܋	܋	܋	܋	܋	܋	܋	܋	܋	܋	܋	܋	
F	܌	܌	܌	܌	܌	܌	܌	܌	܌	܌	܌	܌	܌	܌	܌	

Space
character at
 $0x20 = 32_{(10)}$

Lowercase 'a': $0x97$

Justifies $h_{char}(c) = (int)c - 97$
when lowercase characters are
used!

Other whitespace
chars (\t, LF, ...)

ASCII

Capital 'N' at $0x4E =$
 $79_{(10)}$

I'm the IBM-PC and I
do what I want with
my top bit

OEM standard

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	☺	☻	♥	♦	♣	♠	⌚	⌚	⌚	⌚	⌚	⌚	⌚	⌚	⌚	
1	◀	▶	↔	☰	✖	✖	✖	✖	✖	✖	✖	✖	✖	✖	✖	
2	?	"	#	\$	%	&	,	<	>	*	+	,	-	/		
3	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>?	
4	Ѐ	Ѐ	Ѐ	Ѐ	Ѐ	Ѐ	Ѐ	Ѐ	Ѐ	Ѐ	Ӣ	Ӣ	Ӣ	Ӣ	Ӣ	
5	܂	܂	܂	܂	܂	܂	܂	܂	܂	܂	܂	܂	܂	܂	܂	
6	܃	܃	܃	܃	܃	܃	܃	܃	܃	܃	܃	܃	܃	܃	܃	
7	܄	܄	܄	܄	܄	܄	܄	܄	܄	܄	܄	܄	܄	܄	܄	
8	܅	܅	܅	܅	܅	܅	܅	܅	܅	܅	܅	܅	܅	܅	܅	
9	܆	܆	܆	܆	܆	܆	܆	܆	܆	܆	܆	܆	܆	܆	܆	
A	܇	܇	܇	܇	܇	܇	܇	܇	܇	܇	܇	܇	܇	܇	܇	
B	܈	܈	܈	܈	܈	܈	܈	܈	܈	܈	܈	܈	܈	܈	܈	
C	܉	܉	܉	܉	܉	܉	܉	܉	܉	܉	܉	܉	܉	܉	܉	
D	܊	܊	܊	܊	܊	܊	܊	܊	܊	܊	܊	܊	܊	܊	܊	
E	܋	܋	܋	܋	܋	܋	܋	܋	܋	܋	܋	܋	܋	܋	܋	
F	܌	܌	܌	܌	܌	܌	܌	܌	܌	܌	܌	܌	܌	܌	܌	

- This turned out to be a problem...

Space character at
0x20 = 32₍₁₀₎

Lowercase 'a': 0x97

Justifies $h_{char}(c) = (int)c - 97$
when lowercase characters are
used!

Other whitespace
chars (\t, LF, ...)

ASCII

Capital 'N' at 0x4E =
79₍₁₀₎

I'm the IBM-PC and I
do what I want with
my top bit

ANSI Solution: Code pages

- **Axiom: Everybody agrees on the ASCII subset**
 - So, 7 LSBs interpretable in only one way across all computers in the world.

ANSI Solution: Code pages

- **Axiom: Everybody agrees on the ASCII subset**
 - So, 7 LSBs interpretable in only one way across all computers in the world.
- Top bit: Alphabet-specific
 - This led to *different code pages* for different alphabets!
 - E.g Icelandic: Code page #00861, Greek: Code page CP00851
 - So, as long as all Greek computers use ANSI CP00851 (reasonable assumption), they will all be able to transmit strings to one another without problems ☺

Look at these code pages

Icelandic

HEX DIGITS	1ST → 2ND ↓	0-	1-	2-	3-	4-	5-	6-	7-	8-	9-	A-	B-	C-	D-	E-	F-
-0		▶ (SP)	0	@	P	`	p	Ç	É	á	SF140000	SF020000	SF460000	GA010000	SA480000	≡	
-1	☺	◀ !	1 A	Q	a	q	ü	æ	í	SF150000	SF070000	SF470000	LS610000	SA020000	Þ ±		
-2	☻	↕ "	2 B	R	b	r	é	Æ	ó	SF060000	SF160000	SF480000	GG020000	SA530000	≥		
-3	♥	!! #	3 C	S	c	s	â û	ð ú	π	SF110000	SF080000	SF490000	GP010000	SA520000	≤		
-4	♦	¶ \$	4 D	T	d	t	ä ö	Á	Σ	SF090000	SF100000	SF500000	GS020000	SS260000	∫		
-5	♣	§ %	5 E	U	e	u	à þ	Í	σ	SF190000	SF050000	SF510000	GS010000	SS270000	ʃ		
-6	♠	- &	6 F	V	f	v	å û	Ó	μ	SF360000	SF200000	SF520000	GM010000	SA060000	÷		
-7	●	↑ '	7 G	W	g	w	ç	Ý	Ú	SF210000	SF370000	SF530000	GT010000	SA700000	≈		
-8	■	↑ (8 H	X	h	x	ê ý	ł	Φ	SF220000	SF380000	SF540000	GF020000	SM190000	°		
-9	○	↓)	9 I	Y	i	y	ë Ö	¶	Θ	SF100000	SF390000	SF400000	GT620000	SA790000	•		
-A	Ⓐ	→ *	: J	Z	j	z	è Ü	¬	Ω	SF010000	GO320000	SD630000			·		
-B	♂	← + ;	K	[k	{	Đ Ø	½	δ	SF140000	SF250000	SF410000	GD010000	SA800000	✓		
-C	♀	↳ ,	< L	\	l		ð £	¼	∞	SF260000	SF420000	SF570000	SA450000	LN011000	n		
-D	♪	↔ - =	M]	m	}	Þ Ø	ø	í	ϕ	SF270000	SF430000	SF580000	GF010001	ND021000	²		
-E	♪	▲ .	> N	^ n	~	Ä Pts	Þ	« »	ε	SF280000	SF440000	SF590000	GE010000	SM470000			
-F	☀	▼ / ?	O	— o	◊ Å	f	Å	»	□	SF030000	SF450000	SF600000	SA380000	SP300000	(RSP)		

Greek

HEX DIGITS	1ST → 2ND ↓	0-	1-	2-	3-	4-	5-	6-	7-	8-	9-	A-	B-	C-	D-	E-	F-
-0		▶ (SP)	0	@	P	`	p	Α	Ρ	ι	SF140000	SF020000	SF460000	GO310000	GO720000	Ω	
-1	☺	◀ !	1 A	Q	a	q	ü	æ	ί	SF150000	SF070000	SF470000	GR020000	GI010000	ά ±		
-2	☻	↕ "	2 B	R	b	r	é	Æ	ό	SF060000	SF160000	SF480000	GG020000	GA530000	≥		
-3	♥	!! #	3 C	S	c	s	â ô	ú	π	SF110000	SF080000	SF490000	GP010000	SA520000	≤		
-4	♦	¶ \$	4 D	T	d	t	ä ö	Á	Σ	SF090000	SF100000	SF500000	GS020000	GE710000	Ϊ		
-5	♣	§ %	5 E	U	e	u	à þ	Í	σ	SF190000	SF050000	SF510000	GS010000	SS260000	Ϋ		
-6	♠	- &	6 F	V	f	v	å û	Ó	μ	SF360000	SF200000	SF520000	GM010000	SA060000	÷		
-7	●	↑ '	7 G	W	g	w	ç	Ý	Ú	SF210000	SF370000	SF530000	GT010000	SA170000	≈		
-8	■	↑ (8 H	X	h	x	ê ý	ł	Φ	SF220000	SF380000	SF540000	GF020000	SM190000	°		
-9	○	↓)	9 I	Y	i	y	ë Ö	¶	Θ	SF100000	SF390000	SF400000	GT620000	SA790000	•		
-A	Ⓐ	→ *	: J	Z	j	z	è Ü	¬	Ω	SF010000	GO320000	SD630000			·		
-B	♂	← + ;	K	[k	{	Đ Ø	½	δ	SF140000	SF250000	SF410000	GD010000	SA800000	✓		
-C	♀	↳ ,	< L	\	l		ð £	¼	∞	SF260000	SF420000	SF570000	SA450000	LN011000	n		
-D	♪	↔ - =	M]	m	}	Þ Ø	ø	í	ϕ	SF270000	SF430000	SF580000	GF010001	ND021000	²		
-E	♪	▲ .	> N	^ n	~	Ä Pts	Þ	« »	ε	SF280000	SF440000	SF590000	GE010000	SM470000			
-F	☀	▼ / ?	O	— o	◊ Å	f	Å	»	□	SF030000	SF450000	SF600000	SA380000	SP300000	(RSP)		

Look.

Problems with code page scheme

1. Advent of the web and e-mail: browsers and mail programs would need to store all code books in order to be able to decode messages from different countries correctly!
 - Back in the day, when memory was expensive, this required a lot of space.
 - E.g Jason wants to use his Mac, manufactured in China, according to American specifications, to e-mail his mother, who is currently in China, using a laptop purchased in Greece. Oh, and he types the e-mail in Greek!
2. For some languages, e.g Mandarin, 8 bits are simply not enough!

The Unicode standard

- Defined and maintained by [the Unicode Consortium](#).
- An effort to develop a “**mega-code page**” to satisfy all users in the world.
- Every character gets mapped to a single “**code point**” which is just a number in Hex followed by the string “U+” (stands for Unicode).
 - E.g U+0388= λ (*Greek lowercase lambda*), U+0153=œ (*latin ‘oe’*)
 - “**Hello**” = **U+0048 U+0065 U+006C U+006C U+006F**.

The Unicode standard

- Defined and maintained by [the Unicode Consortium](#).
- An effort to develop a “**mega-code page**” to satisfy all users in the world.
- Every character gets mapped to a single “**code point**” which is just a number in Hex followed by the string “U+” (stands for Unicode).
 - E.g U+0388= λ (*Greek lowercase lambda*), U+0153=œ (*latin ‘oe’*)
 - “**Hello**” = **U+0048 U+0065 U+006C U+006C U+006F**.

The Unicode standard

- “Hello” = U+0048 U+0065 U+006C U+006C U+006F.
- **This hex number does not necessarily represent the number of bits required to store the character in memory!**
 - E.g U+10FFFF is the largest Unicode codepoint and clearly it can not fit into either one or two bytes...
- How the code point itself is mapped to memory (the code point’s *encoding*) is left to the destination computer ☺

Implementation of encodings

- **STEP 0: EVERYBODY AGREES ON USING UNICODE. EVERYONE.**

Implementation of encodings

- **STEP 0: EVERYBODY AGREES ON USING UNICODE. EVERYONE.**
- Step 1: Select an encoding.

Implementation of encodings

- **STEP 0: EVERYBODY AGREES ON USING UNICODE. EVERYONE.**
- Step 1: Select an encoding.
- Many different encodings of Unicode exist, and are pre-installed in all major OSs.
 - UCS2 and its successor UTF-16 stores the codepoint U+xxxx using exactly xxxx, so in 16 bits.
 - UTF8 stores every codepoint in MULTIPLES OF 1 BYTE. (*variable-length encoding*)
 - Codepoints from U+0x0000 to U+0x7E ($127_{(10)}$) are stored **within an entire byte of their own (so the MSB is effectively ignored!)**.
 - Codepoints above $127_{(10)}$ use as many bytes as they might need!

Implementation of encodings

- **STEP 0: EVERYBODY AGREES ON USING UNICODE. EVERYONE.**
- Step 1: Select an encoding.
- Many different encodings of Unicode exist, and are pre-installed in all major OSs.
 - UCS2 and its successor UTF-16 stores the codepoint U+xxxx using exactly xxxx, so in 16 bits.
 - UTF8 stores every codepoint in MULTIPLES OF 1 BYTE. (*variable-length encoding*)
 - Codepoints from U+0x0000 to U+0x7E ($127_{(10)}$) are stored **within an entire byte of their own (so the MSB is effectively ignored!).**
 - Codepoints above $127_{(10)}$ use as many bytes as they might need!
- Step 2: The sender needs to tell the recipient **what kind of Unicode encoding it used to build the text**
 - E-mail headers are a good example
 - In HTML, we can use <meta> tags

Quiz

Quiz

- UTF-8 is racist

Quiz

- UTF-8 is racist

YES

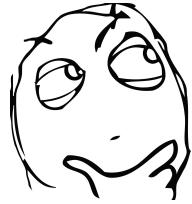
NO

WHAT?

UTF-8 is kind of bigoted ☹

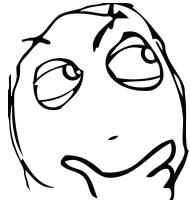
UTF-8 is kind of bigoted ☹

- Or is it?



UTF-8 is kind of bigoted 😞

- Or is it?



- Vast majority of pairs of receivers – senders in the world will send and read characters in English (so, from the set of the first 128 characters that UTF-8 represents in a single byte)
 - This is exactly ASCII! 😊

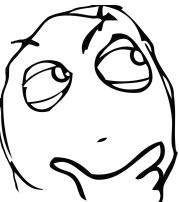
UTF-8 is kind of bigoted 😞

- Or is it? 
- Vast majority of pairs of receivers – senders in the world will send and read characters in English (so, from the set of the first 128 characters that UTF-8 represents in a single byte)
 - This is exactly ASCII! 😊
- As a result, the vast majority of messages can be encoded in one byte per character.
 - We can't hope to do any better when encoding a single character! 😊

UTF-8 is kind of bigoted 😞

- Or is it? 
- Vast majority of pairs of receivers – senders in the world will send and read characters in English (so, from the set of the first 128 characters that UTF-8 represents in a single byte)
 - This is exactly ASCII! 😊
- As a result, the vast majority of messages can be encoded in one byte per character.
 - We can't hope to do any better when encoding a single character! 😊
- Non-English speakers or people who want to use rare characters are disadvantaged: additional storage cost has to be paid!

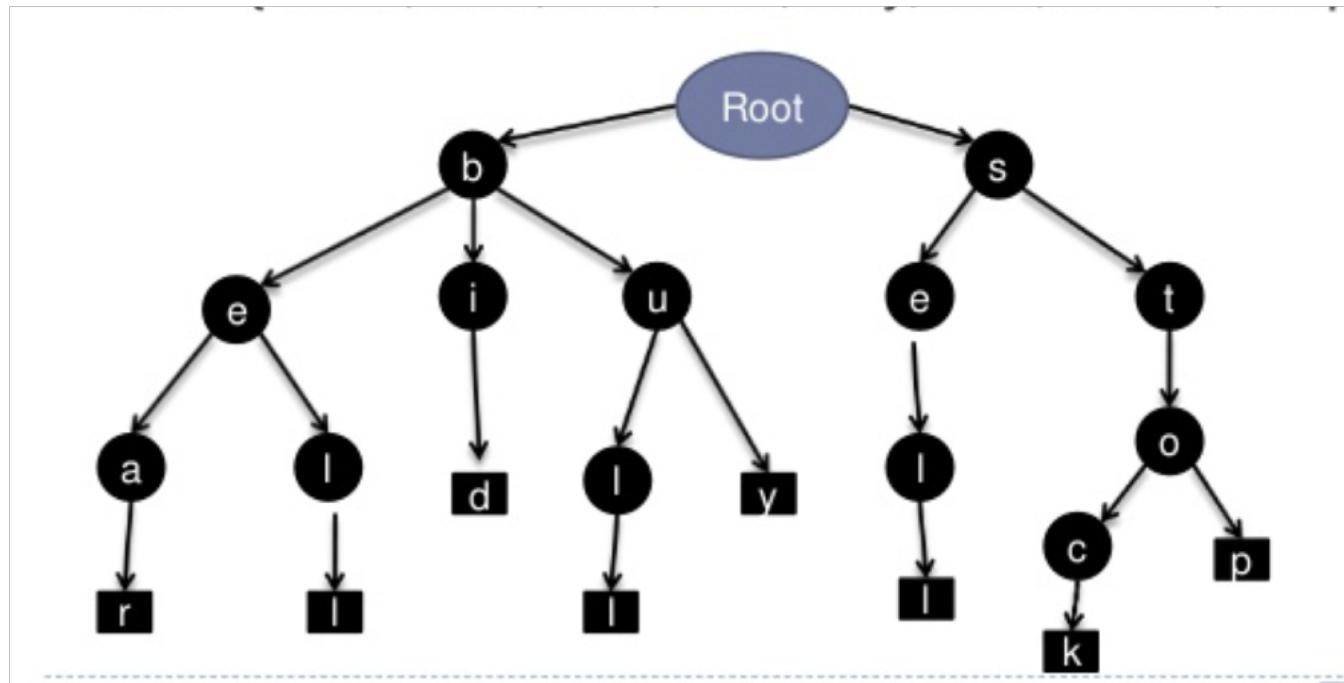
UTF-8 is kind of bigoted 😞

- Or is it? 
- Vast majority of pairs of receivers – senders in the world will send and read characters in English (so, from the set of the first 128 characters that UTF-8 represents in a single byte)
 - This is exactly ASCII! 😊
- As a result, the vast majority of messages can be encoded in one byte per character.
 - We can't hope to do any better when encoding a single character! 😊
- Non-English speakers or people who want to use rare characters are disadvantaged: additional storage cost has to be paid!
 - **Experimentation:** try to load a webpage in Mandarin, making sure that the HTML source reads `<meta http-equiv="utf8"` right after the `<html>` tag... It should take longer than loading the Washington Post, for instance.

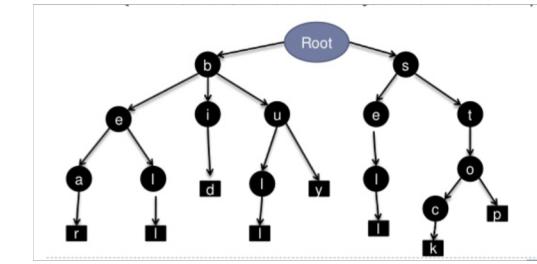
ASCII is good enough for us!

- For our study of strings, **128 characters will be enough.**
- So we limit ourselves to ASCII (technically, 7 of the 8 bits contained in the byte reserved via our assumed UTF-8 encoding).
- ‘A’ in ASCII: $0x41 = 65_{(10)}$
- ‘a’ in ASCII: $0x61 = 97_{(10)}$
- Since $65 + 24 = 89 < 97$, some characters in between lowercase and uppercase alphabet...
- Subtracting 65 from all of our code points allows us to treat ‘A’ at 0...

Tries

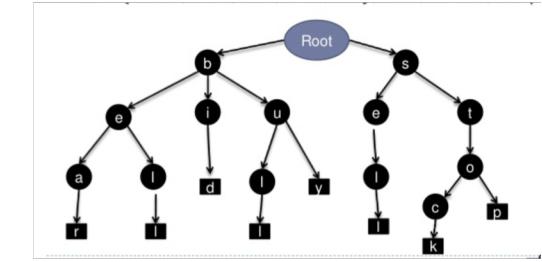


Tries



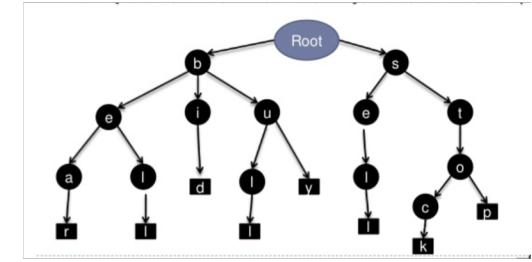
- Pronounced exactly like the verb!

Tries



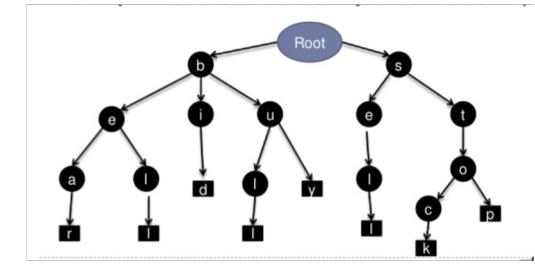
- Pronounced **exactly like the verb!**
 - E.g “José **tries** to catch the ball”
 - This can be sometimes perplexing, since they are essentially *trees*... don’t let this confuse you.

Tries



- Pronounced **exactly like the verb!**
 - E.g “José **tries** to catch the ball”
 - This can be sometimes perplexing, since they are essentially *trees*... don’t let this confuse you.
- A trie is a tree-like data structure where **paths store strings**, and branches are taken based on the relevant **character** of the string that is searched for (or inserted / deleted).
 - An additional bit field tells us **if we have reached an actual string that was stored in the trie.**

Tries



- Pronounced **exactly like the verb!**
 - E.g “José **tries** to catch the ball”
 - This can be sometimes perplexing, since they are essentially *trees*... don’t let this confuse you.
- A trie is a tree-like data structure where **paths store strings**, and **branches are taken based on the relevant character of the string** that is searched for (or inserted / deleted).
 - An additional bit field tells us **if we have reached an actual string that was stored in the trie**.
- Invented in late 1950s by René de la Briandais, esteemed French person.

Structure of a trie node

- We must always specify the alphabet Σ that we sample characters from!
 - Common choices: ASCII, UNICODE.
- In this course, we will be using ASCII or small subsets thereof.
- Every node will hold a buffer with $|\Sigma|$ –many cells! 
- Question: For UNICODE ($\approx 65K$ usable symbols), how big will those buffers be in bytes?

Structure of a trie node

- We must always specify the alphabet Σ that we sample characters from!
 - Common choices: ASCII, UNICODE.

Structure of a trie node

- We must always specify the alphabet Σ that we sample characters from!
 - Common choices: ASCII, UNICODE.
- In this course, we will be using ASCII or small subsets thereof.

Structure of a trie node

- We must always specify the alphabet Σ that we sample characters from!
 - Common choices: ASCII, UNICODE.
- In this course, we will be using ASCII or small subsets thereof.
- Every node will hold a buffer with $|\Sigma|$ –many cells!



Structure of a trie node

- We must always specify the alphabet Σ that we sample characters from!
 - Common choices: ASCII, UNICODE.
- In this course, we will be using ASCII or small subsets thereof.
- Every node will hold a buffer with $|\Sigma|$ –many cells! 
- Question: For UNICODE ($\approx 65K$ usable symbols), how big will those buffers be in bytes?

128

65K

>65K

<65K

Structure of a trie node

- We must always specify the alphabet Σ that we sample characters from!
 - Common choices: ASCII, UNICODE.
- In this course, we will be using ASCII or small subsets thereof.
- Every node will hold a buffer with $|\Sigma|$ –many cells! 
- Question: For UNICODE ($\approx 65K$ usable symbols), how big will those buffers be in bytes?

128

65K

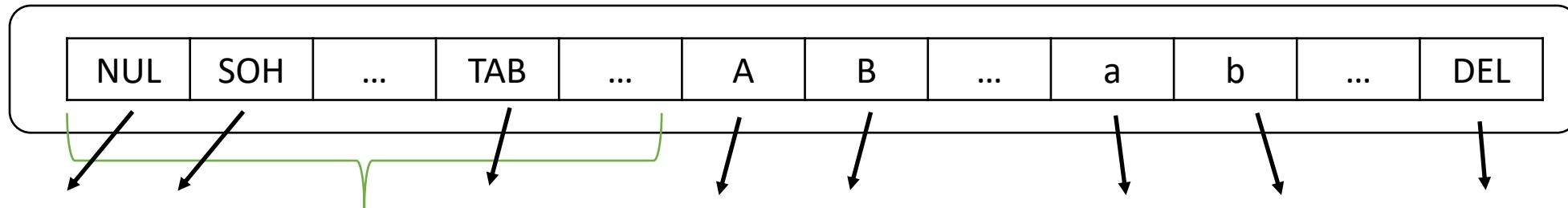
>65K

<65K

Only the first 128 characters are stored in a single byte each!

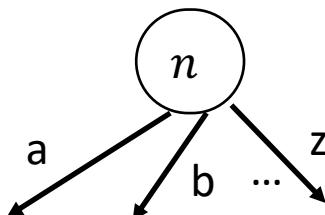
Structure of a trie node

- This is an example of what a trie node would look like if alphabet is ASCII:



Non-printable characters, like system bell, paper feed, carriage return....

- To make the art (?) tractable, we will be shortening the nodes to only have outgoing links that correspond to actual characters that are **dynamically used** in the trie



Structure of a trie node

- It turns out that in a trie node, we also need some information that tells us

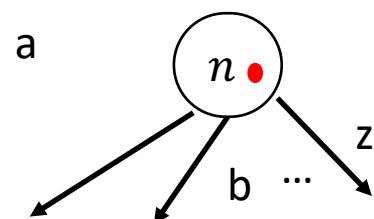
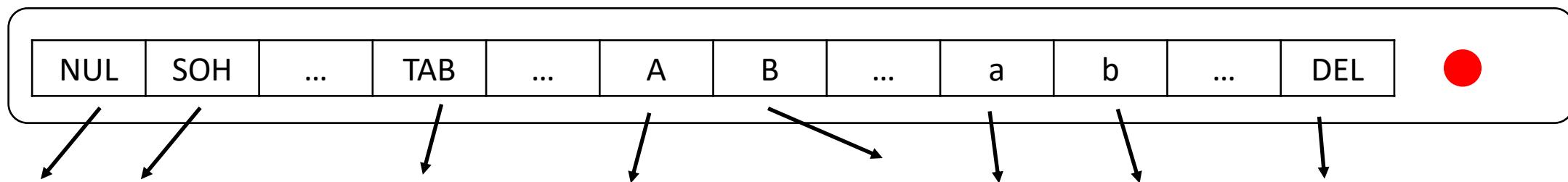
Does the path from the root of the trie correspond to an actual string stored in it?

Structure of a trie node

- It turns out that in a trie node, we also need some information that tells us

Does the path from the root of the trie correspond to an actual string stored in it?

- This corresponds to a single bit! “Art” convention: Red for set, blue for unset



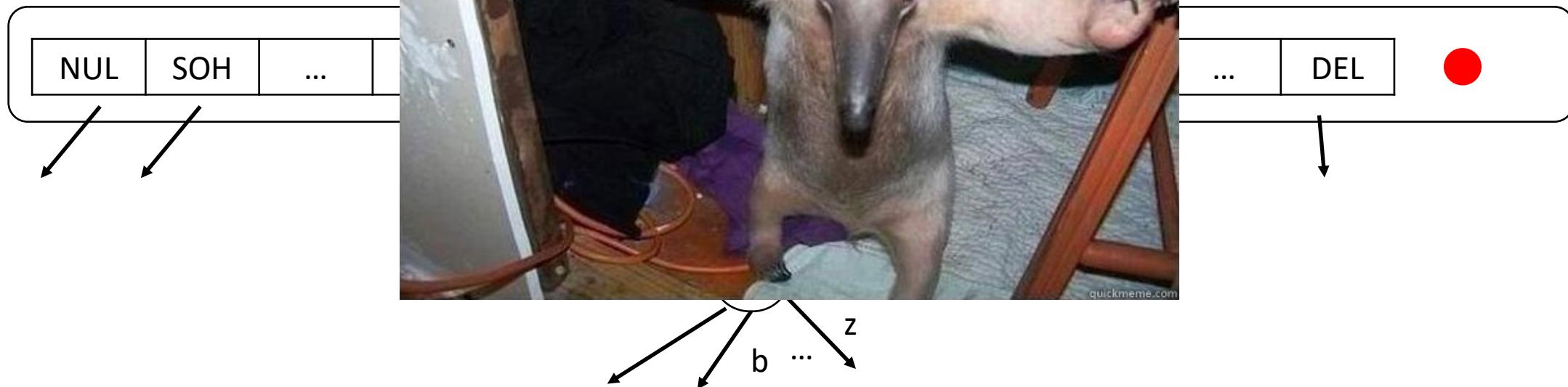
Structure of a trie node

- It turns out that in a trie node, we also need some information that tells us

Does the path from the root to this node represent a valid string?

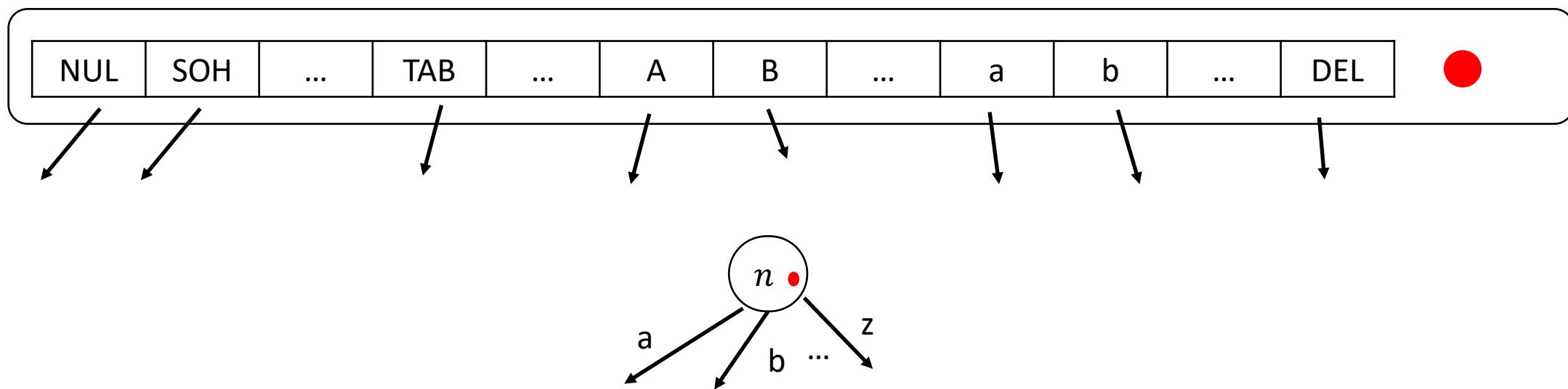
'string stored in it?

- This corresponds to a



Not so fast!

- Don't I need **2 buffers**, to store both characters **and** pointers?

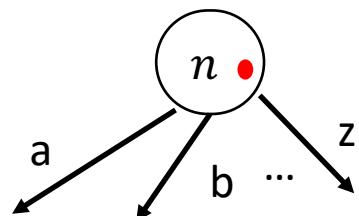
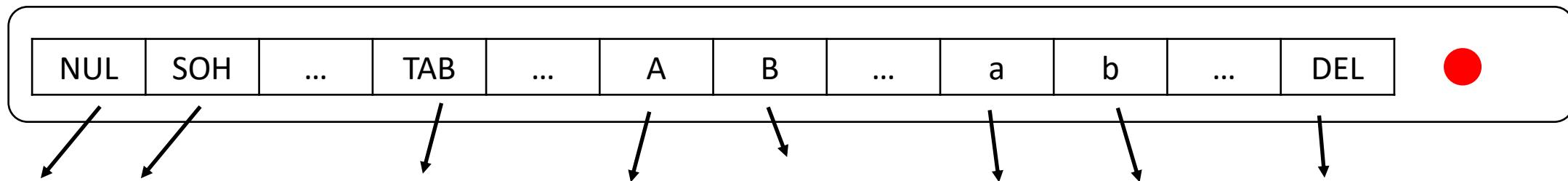


Not so fast!

- Don't I need **2 buffers**, to store both characters **and** pointers?

Yes

No



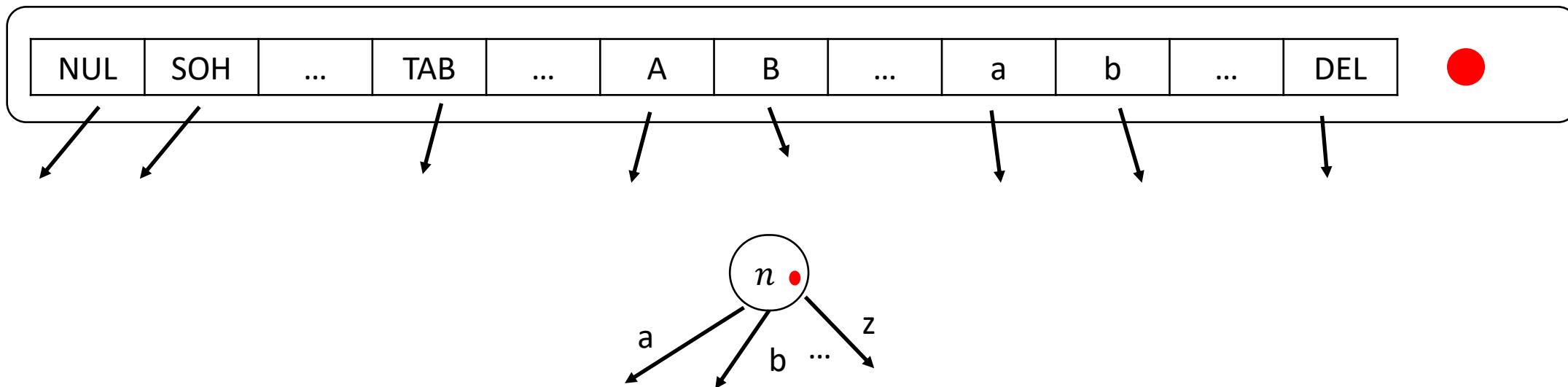
Not so fast!

- Don't I need **2 buffers**, to store both characters **and** pointers?

Yes

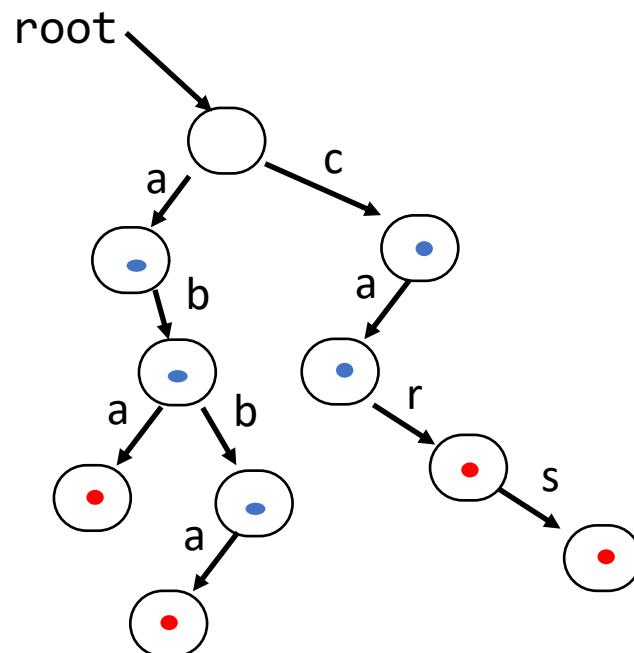
No

- Every character gets mapped to a number through ASCII! (e.g `(int)`A' = 65`)
- Can absolutely do something like `buf['b']`! ☺



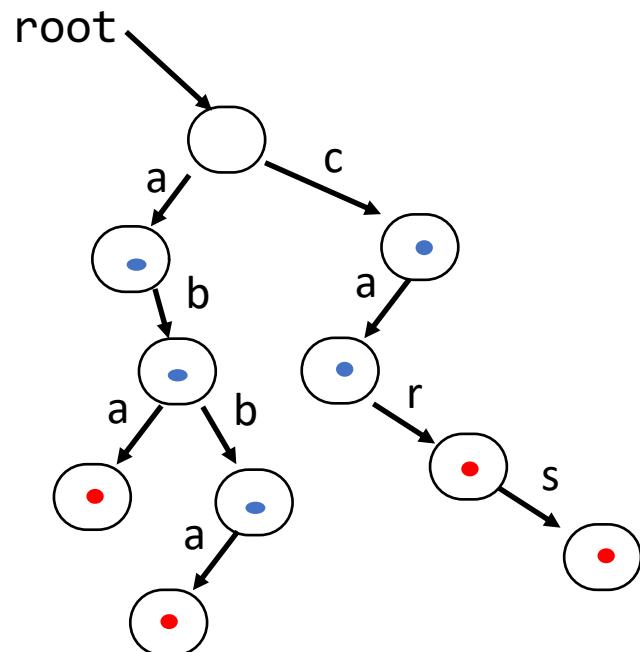
Example

- Assume alphabet to be ASCII: 128 characters.
- The following trie stores the strings “abba”, “aba”, “car”, “cars”:



Example

- Assume alphabet to be ASCII: 128 characters.
- The following trie stores the strings “abba”, “aba”, “car”, “cars”:

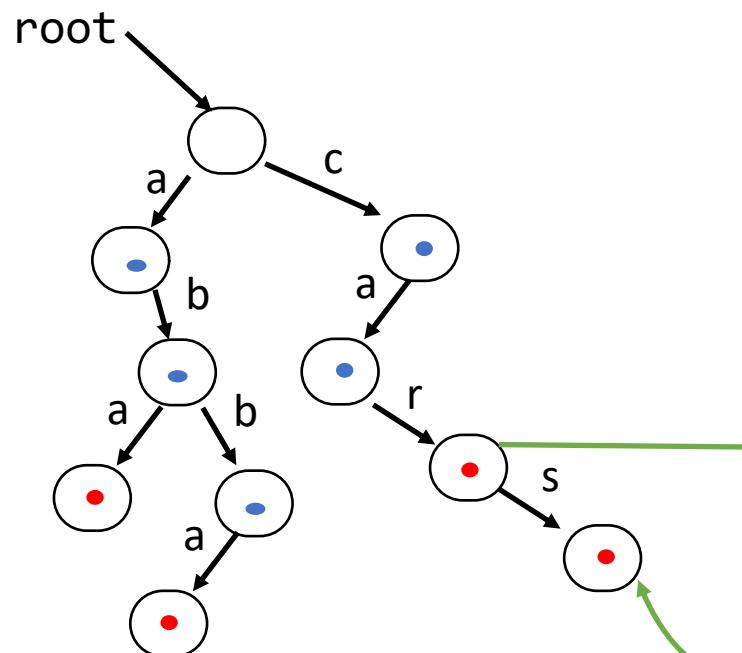


Do we perhaps see the utility of the “**end-of-stored string**” bit now? 😊

Example



- Assume alphabet to be ASCII: 128 characters.
- The following trie stores the strings “abba”, “aba”, “car”, “cars”:



Do we perhaps see the utility of the “**end-of-stored string**” bit now? 😊

Note that encountering a set bit while traversing the tree is **not a sufficient reason** to stop your search for more keys along that path!

Inserting a key

- Quiz: Inserting a new key into a trie always involves the allocation of new nodes.

Yes
(why?)

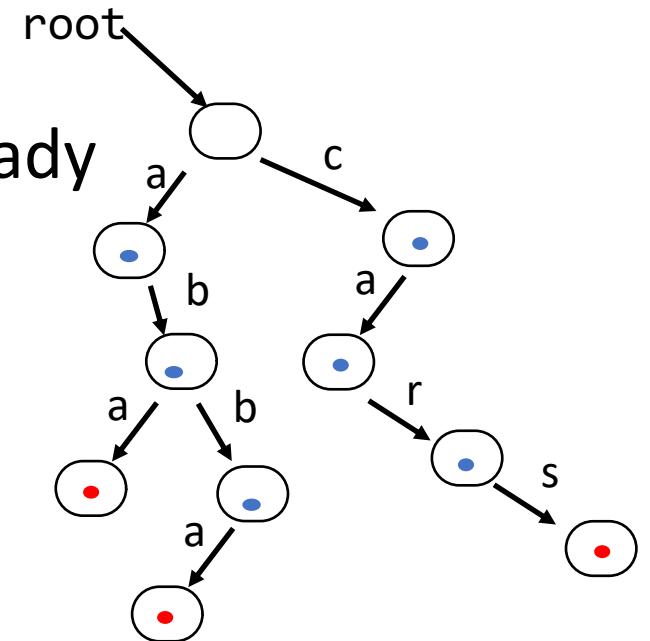
No
(give me an example)

Inserting a key

- Quiz: Inserting a new key into a trie always involves the allocation of **new nodes**.



- In our familiar trie, suppose that “car” was not already stored.



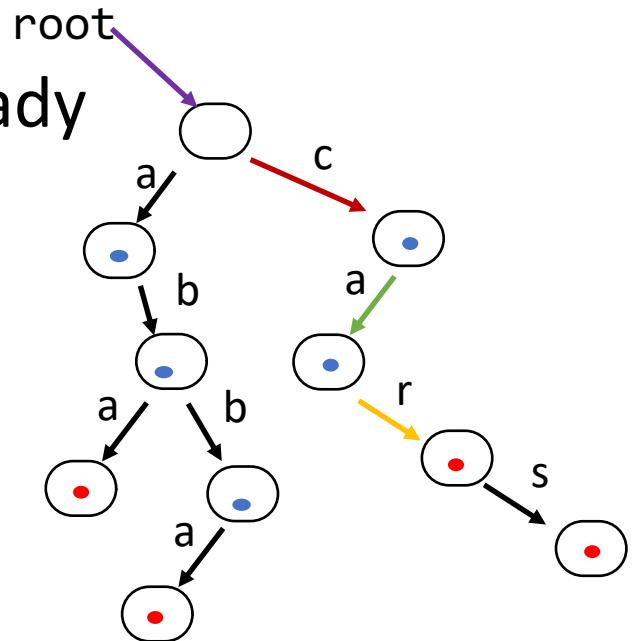
Inserting a key

- Quiz: Inserting a new key into a trie always involves the allocation of new nodes.

Yes
(why?)

No
(give me an example)

- In our familiar trie, suppose that “car” was not already stored.
- Then, we traverse the trie following the links for ‘c’, ‘a’ and ‘r’ and we set the relevant bit!

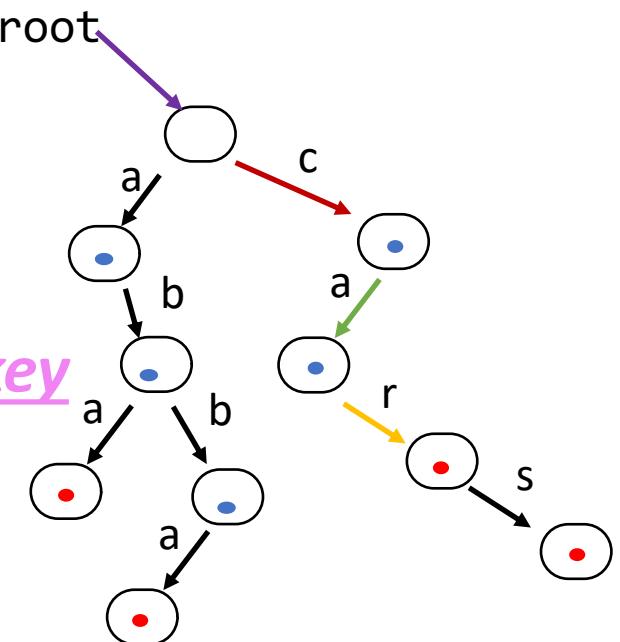


Inserting a key

- Quiz: Inserting a new key into a trie always involves the allocation of **new nodes**.



- In our familiar trie, suppose that “car” was not already stored.
- Then, we traverse the trie following the links for ‘c’, ‘a’ and ‘r’ and **we set the relevant bit!**
- Generally: inserting a key that is a **prefix of an existing key** never leads to a heap visit!

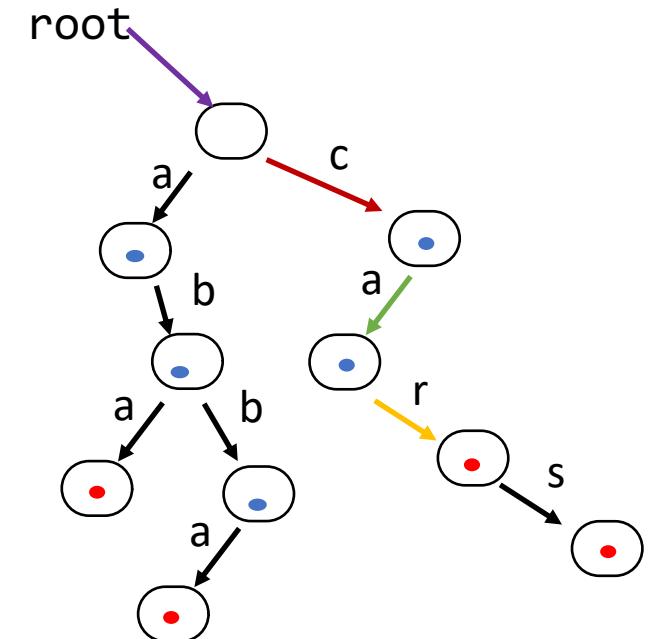


Inserting a key

- Quiz: Inserting a new key into a trie always involves the allocation of **new nodes**.



- In our familiar trie, suppose that “car” was not already stored.
- Then, we traverse the trie following the links for ‘c’, ‘a’ and ‘r’ and **we set the relevant bit!**
- Generally: inserting a key that is a **prefix of an existing key** never leads to a heap visit!
 - With ASCII and Unicode-sized buffers, that’s not too bad of an idea....



Trie insertion: exercise

```
public class AsciiTrie { // No need to make it a generic ; tries store strings!
```

```
    private class Node {
```

```
        private Node[] next = new Node[128];
```

```
        byte eosFlag; // Wish I could manipulate a single bit...
```

```
}
```

```
    private Node root;
```

```
    public void insert(String key){ // More general solution: CharSequences instead of Strings  
        root = insert(root, key, 0); // Could also make it a method of Node
```

```
}
```

```
    private Node insert(Node root, String key, int currStringIndex){
```

```
        /* YOU fill this in NOW!
```

```
        * Hints:
```

```
        *
```

```
index of the string      * (1) Strings have a method with signature charAt(int index) that returns the char at position  
                           *           (index is 0-based).
```

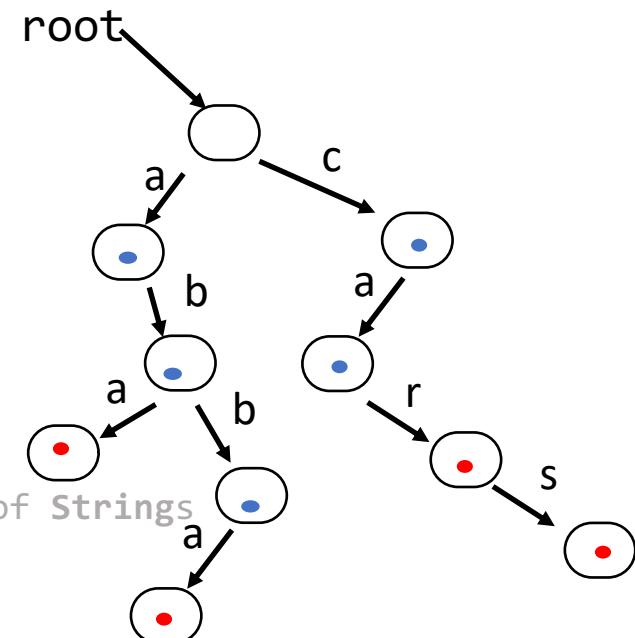
```
        * (2) Think: When do I complete my insertion?
```

```
        * (3) What if I need to allocate a new node?
```

```
        * (4) What if I don't?
```

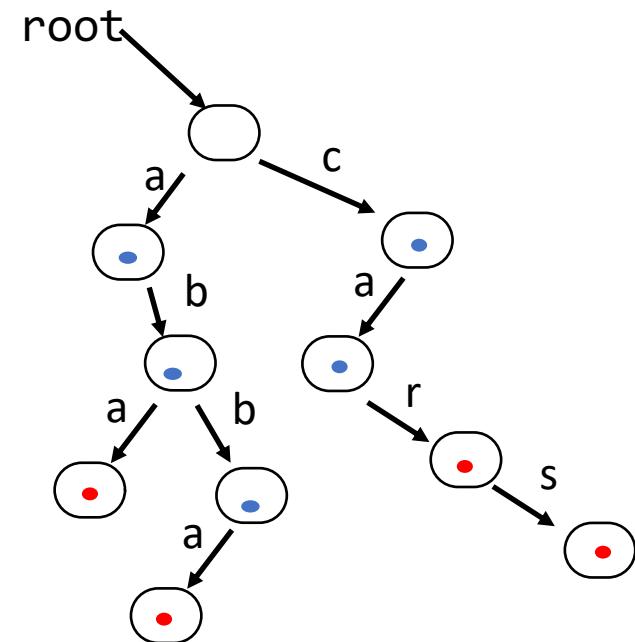
```
        */
```

```
}
```



Trie insertion: exercise

```
public class AsciiTrie { // No need to make it a generic ; tries store strings!
    private class Node {
        private Node[] next = new Node[128];
        byte eosFlag;
    }
    private Node root;
    public void insert(String key){ // More general solution: CharSequences instead of Strings
        root = insert(root, key, 0); // Could also make it a method of Node
    }
    private Node insert(Node node, String key, int currStringIndex){
        assert key != null && key.size() > 0, "Key can't be null or empty.";
        if(node == null) // key is not a prefix of any already stored key...
            node = new Node(); // Node class defn allocates memory immediately, no need for explicit constructor defn.
        if(currStringIndex == key.size())
            node.eosFlag = (byte)1;
        else
            node.next[key.charAt(currStringIndex)] = insert(node.next[key.charAt(currStringIndex)], key, currStringIndex + 1);
    }
    return node;
}
```

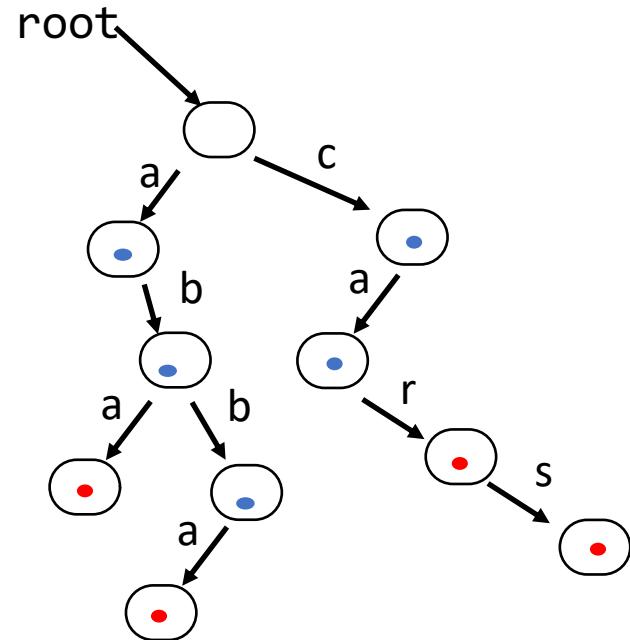


This could
work! ☺



Trie insertion: exercise

```
public class AsciiTrie { // No need to make it a generic ; tries store strings!
    private class Node {
        private Node[] next = new Node[128];
        byte eosFlag;
    }
    private Node root;
    public void insert(String key){ // More general solution: CharSequences instead of Strings
        root = insert(root, key, 0); // Could also make it a method of Node
    }
    private Node insert(Node node, String key, int currStringIndex){
        assert key != null && key.size() > 0, "Key can't be null or empty.";
        if(node == null)
            node = new Node();
        if(currStringIndex == key.size())
            node.eosFlag = (byte)1;
        else
            node.next[key.charAt(currStringIndex)] = insert(node.next[key.charAt(currStringIndex)], key,
                currStringIndex + 1);
    }
}
```



- **key** is not really stored anywhere...
- Can we perhaps get rid of it?

Yes
(how?)

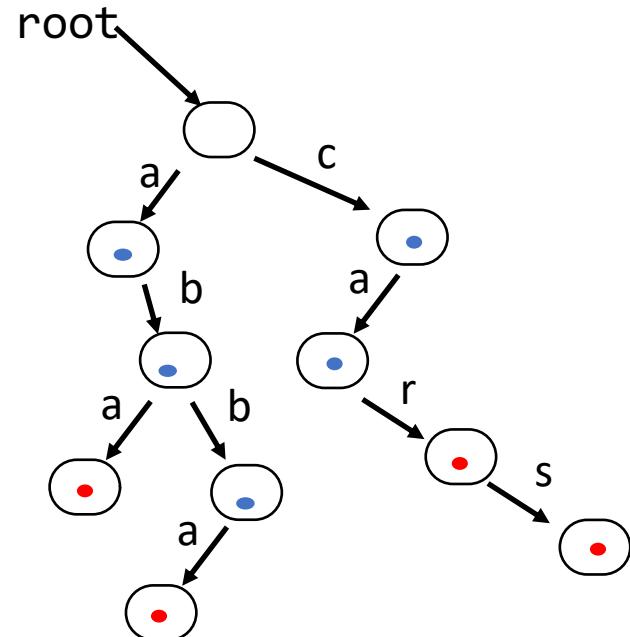
No
(why?)

Trie insertion: exercise

```
public class AsciiTrie { // No need to make it a generic ; tries store strings!
    private class Node {
        private Node[] next = new Node[128];
        byte eosFlag;
    }
    private Node root;
    public void insert(String key){ // More general solution: CharSequences instead of Strings
        root = insert(root, key, 0); // Could also make it a method of Node
    }
    private Node insert(Node node, String key, int currStringIndex){
        assert key != null && key.size() > 0, "Key can't be null or empty.";
        if(node == null){ // key is not a prefix of any already stored key...
            node = new Node();
            if(currStringIndex == key.size())
                node.eosFlag = (byte)1;
            else
                node.next[key.charAt(currStringIndex)] = insert(node.next[key.charAt(currStringIndex)], key, currStringIndex + 1);
        }
        return node;
    }
}
```

Just because **the reference** isn't stored doesn't mean that:

1. It isn't useful for the insertion algorithm
2. There is no **implicit** storage for the string in the trie ([links and node bits!](#))

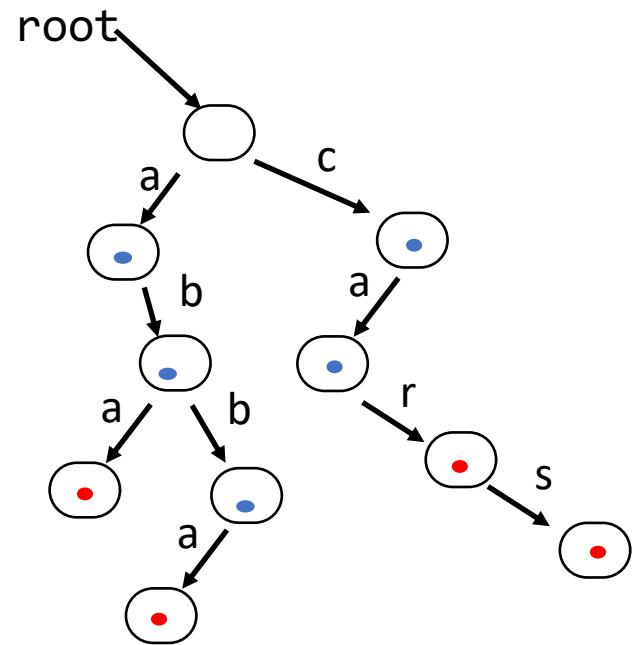


Yes
(how?)

No
(why?)

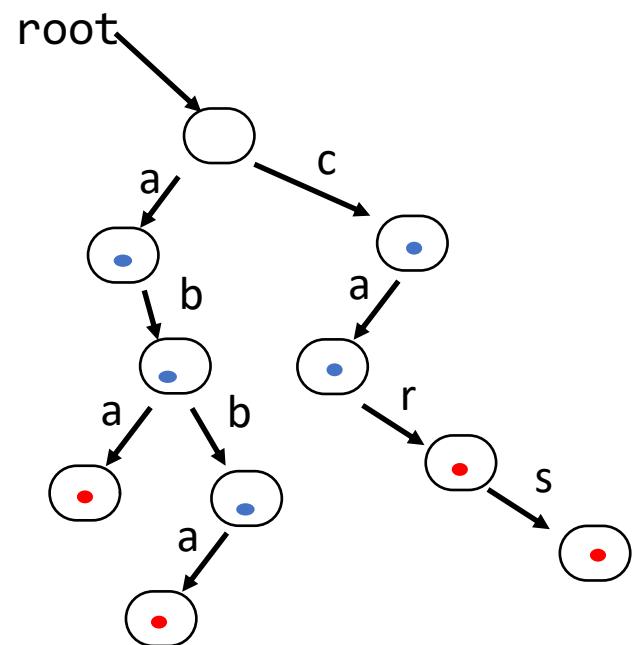
Iterative insertion!

```
public class AsciiTrie { // No need to make it a generic ; tries store strings!
    private class Node {
        private Node[] next = new Node[128];
        byte eosFlag;
    }
    private Node root;
    public void insert(String key){
        /* Fill this in with an iterative implementation
         * of trie insertion!
        */
    }
}
```



Iterative insertion!

```
public class AsciiTrie { // No need to make it a generic ; tries store strings!
    private class Node {
        private Node[] next = new Node[128];
        byte eosFlag;
    }
    private Node root;
    public void insert(String key){
        assert key != null && key.size() > 0, "Key can't be empty.";
        int currIndex = 0;
        Node curr = root, prev = null;
        while(currIndex < key.size()){
            char charToFollow = key.charAt(currIndex);
            if(curr == null)
                curr = new Node();
            if(prev != null)
                prev.next[charToFollow] = curr;
            currIndex++;
            prev = curr;
            curr = curr.next[charToFollow];
        }
        curr.eosFlag = (byte)1;
    }
}
```



Deleting a key

- Quite simple: traverse the trie and, if you find the key, **set the bit to blue**.
- Nothing to do if the bit is already blue (key not actually in DB)

Deleting a key

- Quite simple: traverse the trie and, if you find the key, [set the bit to blue](#).
- Nothing to do if the bit is already blue (key not actually in DB)
- If you encounter a null pointer before you exhaust the key's characters, then you also know that the key is not in DB!

Deleting a key

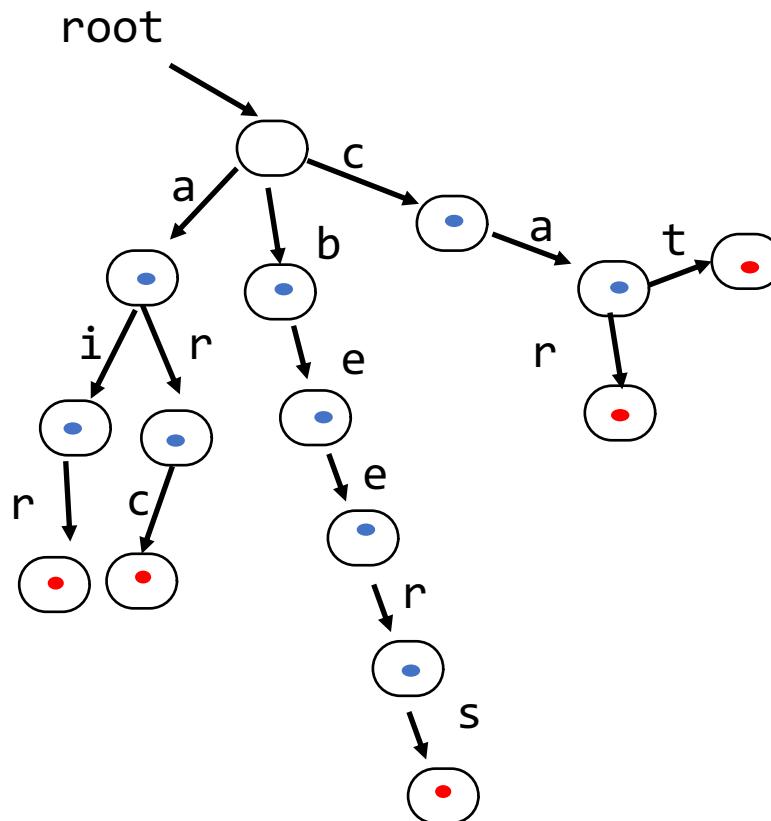
- Quite simple: traverse the trie and, if you find the key, **set the bit to blue**.
- Nothing to do if the bit is already blue (key not actually in DB)
- If you encounter a `null` pointer before you exhaust the key's characters, then you also know that the key is not in DB!
- In those two cases, make a choice between returning `null` or throwing an `Exception`.
 - No matter what you do, **commit to it in your interface!**
 - If you throw an `Exception`, it **must be checked and documented (@throws)**
- This is an example of “soft” deletion.

Deleting a key

- Quite simple: traverse the trie and, if you find the key, **set the bit to blue**.
- Nothing to do if the bit is already blue (key not actually in DB)
- If you encounter a `null` pointer before you exhaust the key's characters, then you also know that the key is not in DB!
- In those two cases, make a choice between returning `null` or throwing an `Exception`.
 - No matter what you do, **commit to it in your interface!**
 - If you throw an `Exception`, it **must be checked and documented (@throws)**
- This is an example of “soft” deletion.

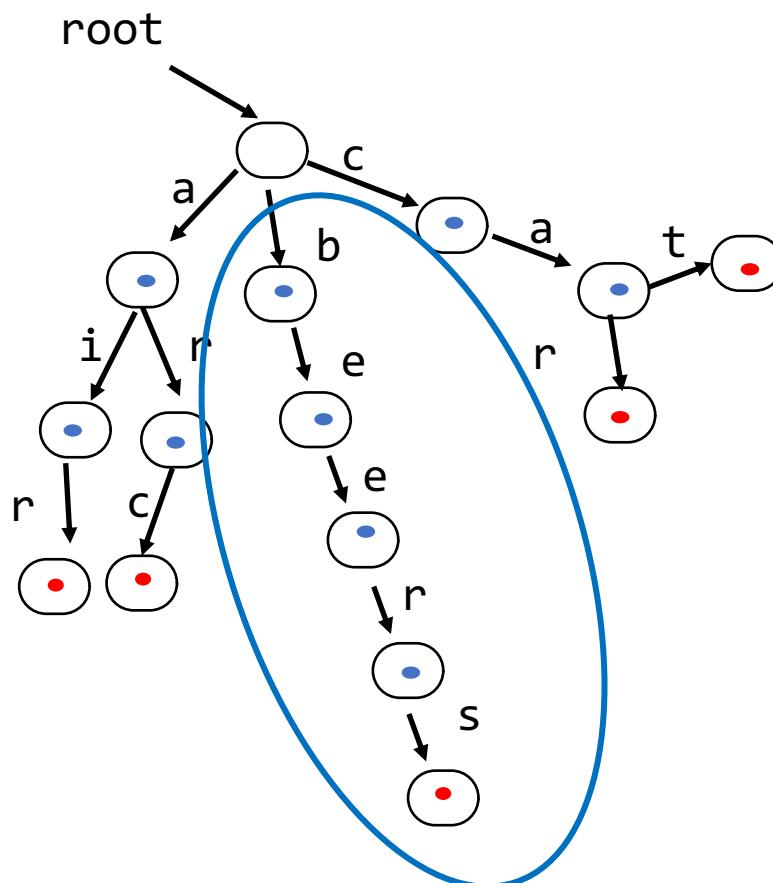
“Soft” vs “hard” deletion in a trie: what would **you** do?

- Suppose that you want to delete ‘‘beers’’ from the following trie.



“Soft” vs “hard” deletion in a trie: what would **you** do?

- Suppose that you want to delete ‘`beers`’ from the following trie.



Since “beers” is the only stored string in this long path from the root (35.7% of total nodes), we could throw away this entire path, e.g by setting `root['b'] = null`

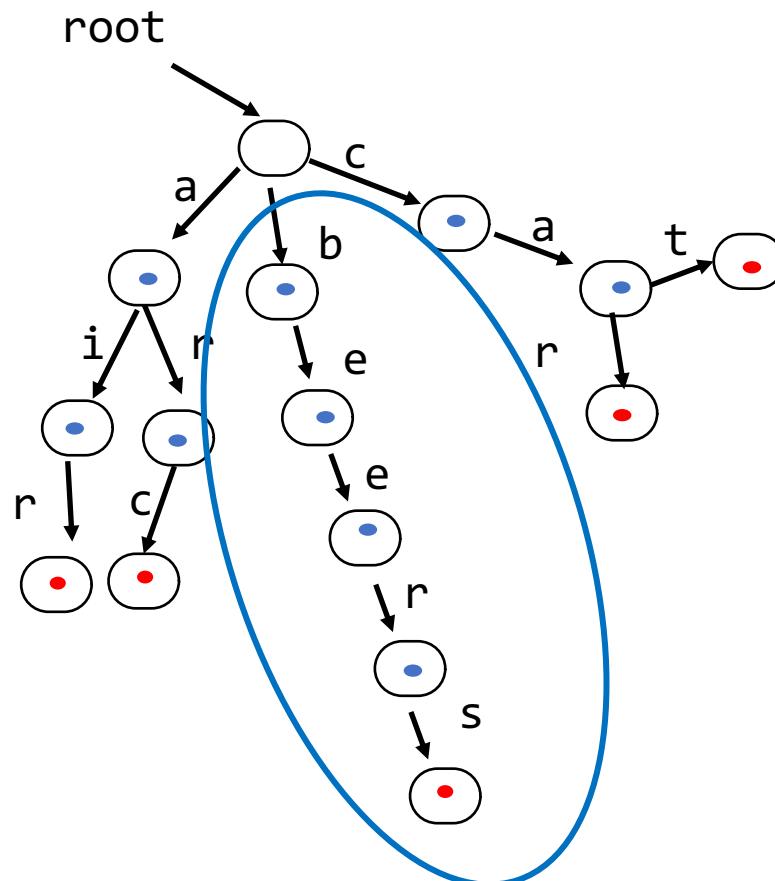
“Soft” vs “hard” deletion in a trie: what would **you** do?

- Suppose that you want to delete ‘`beers`’ from the following trie.

Would **you** do this? (No right or wrong answer here!)

Yes
(why?)

No
(why?)



Since “beers” is the only stored string in this long path from the root (35.7% of total nodes), we could throw away this entire path, e.g by setting `root['b'] = null`

“Soft” vs “hard” deletion in a trie: what would **you** do?

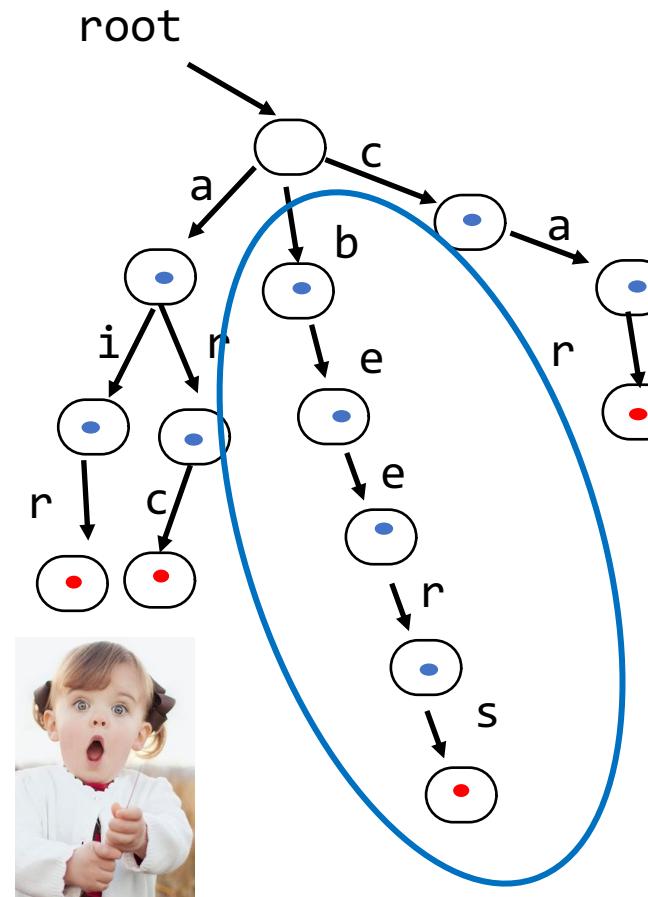
- Suppose that you want to delete ‘‘beers’’ from the following trie.

Would **you** do this? (No right or wrong answer here!)

Yes
(why?)

No
(why?)

Jason **wouldn’t**: strings in a trie, particularly long ones, allow for many **prefixes** to be stored with **either no heap visits or few heap visits!**



E.g consider a future insertion of ‘‘beer’’, ‘‘bee’’, or ‘‘bear!’’

Caveat: If for some reason you know that **in the future you won’t insert any strings that start with a ‘b’**, can delete subtree safely 😊

Complexity of searching for a key in a trie

- Given insertion and deletion, `search(String k)` should be easy for you to implement.

Complexity of searching for a key in a trie

- Given insertion and deletion, `search(String k)` should be easy for you to implement.
- Most important question in tries:** What is the worst-case complexity of searching for a key in a trie with n keys?

$\approx \log_2 n$

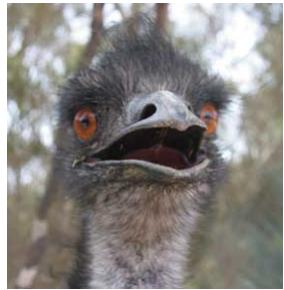
$\approx \log_k n,$
 $k > 2$

$\approx \mathcal{O}(n)$

Something
else
(what?)

Complexity of searching for a key in a trie

- Given insertion and deletion, `search(String k)` should be easy for you to implement.
- Most important question in tries:** What is the worst-case complexity of searching for a key in a trie with n keys?



$\approx \log_2 n$

$\approx \log_k n, k > 2$

$\approx \mathcal{O}(n)$

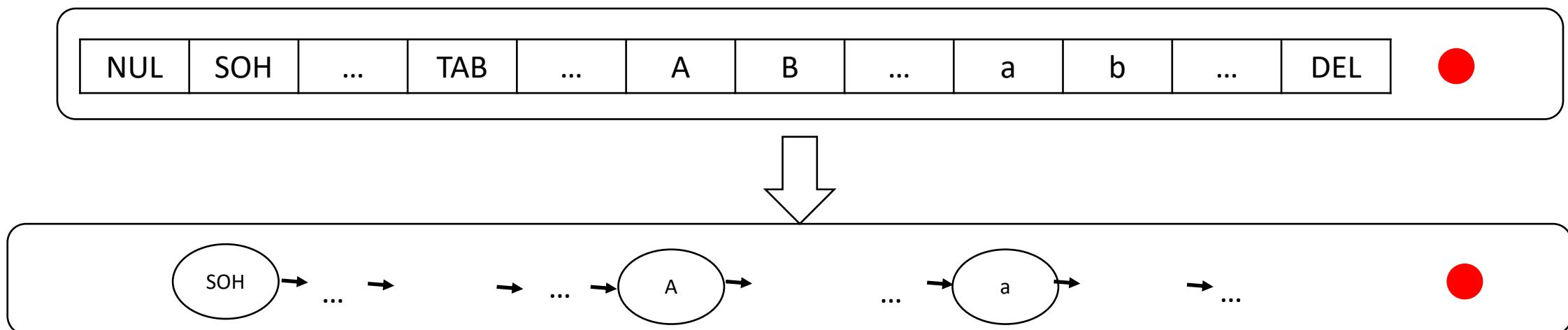
Something
else
(what?)

- Tries: **first data structure we see where search efficiency does not depend on n !!!**

LENGTH OF THE LONGEST EXISTING KEY IN THE TRIE!

De la Briandais Tries

- Obvious problem with tries: they waste a lot of space when characters are not used =/
- Obvious solution: Replace the inner nodes' buffer with a linked list over the characters that have split keys that go through this node!



- Caveat: we lose the $O(1)$ access that `charAt()` was giving us 😞

Cool things we can do with tries

- Quiz: Can we use a trie as a **dictionary**, the way that we've used AVL Trees, Red-Black Trees, etc?

Cool things we can do with tries

- Quiz: Can we use a trie as a **dictionary**, the way that we've used AVL Trees, Red-Black Trees, etc?

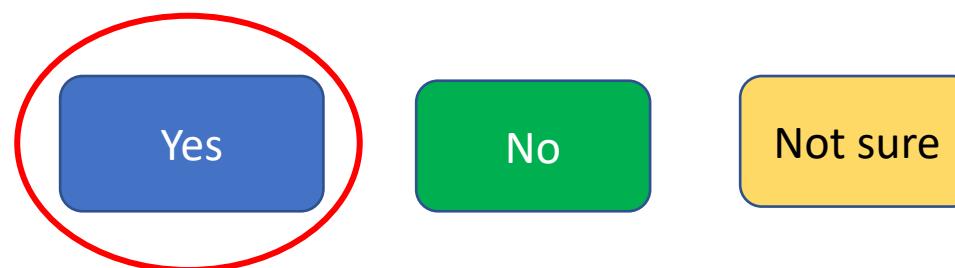
Yes

No

Not sure

Cool things we can do with tries

- Quiz: Can we use a trie as a **dictionary**, the way that we've used AVL Trees, Red-Black Trees, etc?



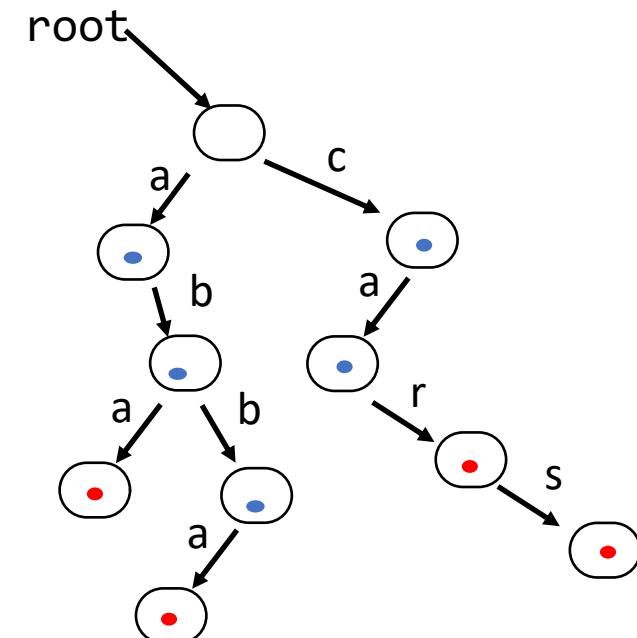
- **Absolutely!**
- They offer **excellent search efficiency** no matter how many strings they hold.
- **Insertion on par with search**
- **Deletion more efficient than B-Trees** (at least for small strings)
 - Plus, as in with BSTs, **all dictionary operations can be implemented iteratively!** 😊

Trie-specific stuff

- Much like with balanced trees, where we could ask for index-specific stuff such as `inorder traversal` and `height`, we can also use tries for non-dictionary related things! ☺
- Practice: In the next slide, fill-in the code of the instance method `longestPrefixOf(String key)`, which returns the longest prefix of key `key` in the trie.

Longest Prefix!

```
public class AsciiTrie {  
    private class Node {  
        private Node[] next = new Node[128];  
        byte eosFlag;  
    }  
    private Node root;  
    public String longestPrefixOf(String key){  
        /* Fill this in!  
         * Hints:  
         * (1) Recall that only nodes with a red bit can be candidates for  
         *     encoding the longest prefix of key!  
         * (2) key.substr(int startIndex, int length) will return a substring of key that  
         *     starts from startIndex and has a length of "length". It might be useful.  
        */  
    }  
}
```



```

public class AsciiTrie {
    private class Node {
        private Node[] next = new Node[128];
        byte eosFlag;
    }
    private Node root;
    public String longestPrefixOf(String key){
        int longestPrefixLength = getMaxPrefixLength(root, key, 0, 0);
        return key.substr(0, longestPrefixLength);
    }
    private int getMaxPrefixLength(Node curr, String key, int currInd, int maxLenRec){
        if(curr == null || currInd == key.length()) return maxLenRec;
        if(curr.eosFlag == (byte)1) maxLenRec = currInd;
        return getMaxPrefixLength(curr.next[key.charAt(currInd)], key, currInd + 1,
            maxLenRec);
    }
}

```

Longest Prefix!

```

public class AsciiTrie {
    private class Node {
        private Node[] next = new Node[128];
        byte eosFlag;
    }
    private Node root;
    public String longestPrefixOf(String key){
        int longestPrefixLength = getMaxPrefixLength(root, key, 0, 0);
        return key.substr(0, longestPrefixLength);
    }
    private int getMaxPrefixLength(Node curr, String key, int currInd, int maxLenRec){
        if(curr == null || currInd == key.length()) return maxLenRec;
        if(curr.eosFlag == (byte)1) maxLenRec = currInd;
        return getMaxPrefixLength(curr.next[key.charAt(currInd)], key, currInd + 1,
                               maxLenRec);
    }
}

```

Longest Prefix!

Metacharacter (glob) search!

```
[jason@MacBook:~/temp ls  
fil.txt      file.txt    file0.txt   file1.txt   file2.txt   myfile.txt  myfile0.txt myfile1.txt  
[jason@MacBook:~/temp ls file*  
file.txt  file0.txt file1.txt file2.txt  
[jason@MacBook:~/temp ls myfile*  
myfile.txt  myfile0.txt myfile1.txt  
jason@MacBook:~/temp
```

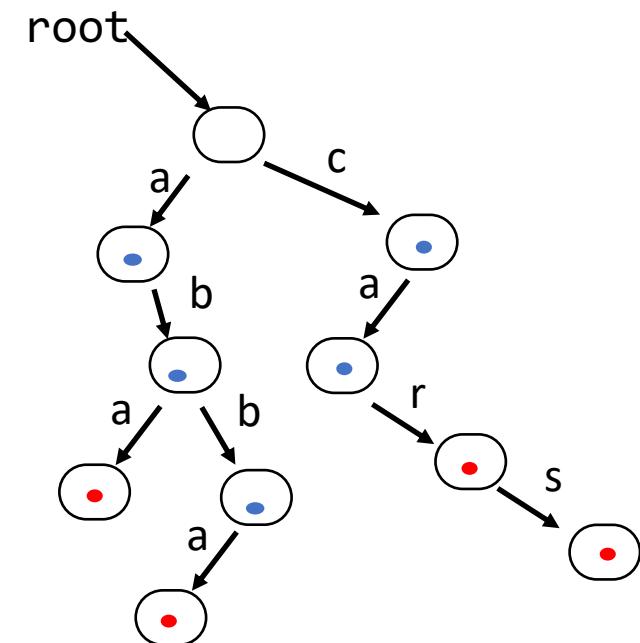
Metacharacter (glob) search!

```
[jason@MacBook:~/temp ls  
fil.txt      file.txt    file0.txt   file1.txt   file2.txt   myfile.txt  myfile0.txt myfile1.txt  
[jason@MacBook:~/temp ls file*  
file.txt  file0.txt file1.txt file2.txt  
[jason@MacBook:~/temp ls myfile*  
myfile.txt  myfile0.txt myfile1.txt  
jason@MacBook:~/temp
```

- The metacharacter '*' will match **every string that begins with “AVL”, and will print them in lexicographic order!**
- The BASH shell stores all file names under the current directory in a trie! 😱
- So let's try to implement **metacharacter parsing** in a trie ☺

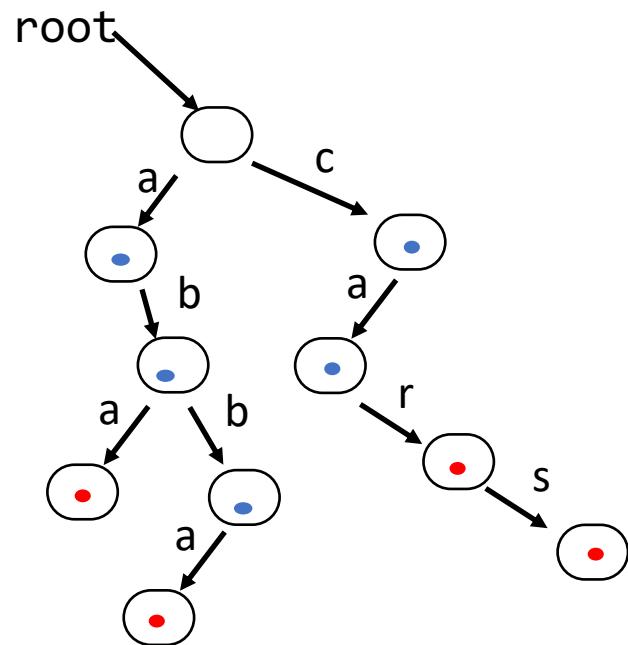
Metacharacter parsing

```
public class AsciiTrie {  
    private class Node {  
        private Node[] next = new Node[128];  
        byte eosFlag;  
    }  
    private Node root;  
    public Collection<String> keysWithPrefix(String prefix){  
        /* Fill the code for this method!  
         *  
         * Hints:  
         *  
         * (1) In the trie to the top right, calling keysWithPrefix("ca") would have the same  
         *     effect as ls ca*. Same thing with ls c*.  
         *  
         * (2) In Java, a Collection is an interface implemented by all kinds of linear  
         *     containers (LinkedList, ArrayList, Queue, Stack, etc)  
         *  
         * (3) Remember: We want Lexicographic ordering!  
        */  
    }  
}
```



Metacharacter parsing

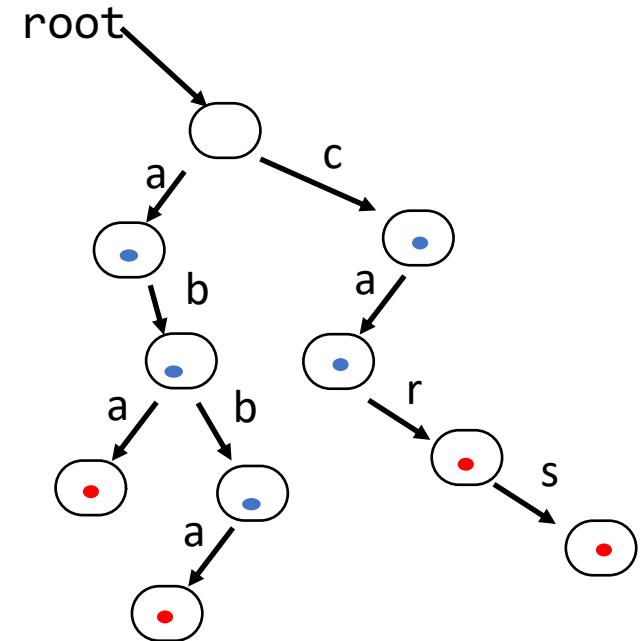
```
public class AsciiTrie {  
    private class Node {  
        private Node[] next = new Node[128];  
        byte eosFlag;  
    }  
    private Node root;  
    public Collection<String> keysWithPrefix(String prefix){  
        Node prefixNode = fetchNode(root, prefix, 0);  
        ArrayList<String> keys = new ArrayList<String>();  
        populate(prefixNode, prefix, keys);  
        return keys;  
    }  
    private void populate(Node curr, String prefix, ArrayList<String> keys){  
        if(curr == null) return;  
        if(curr.eosFlag == (byte)1) keys.add(prefix);  
        for(int i = 0; i < curr.next.length; i++)  
            populate(curr.next[i], prefix + (char)i, keys);  
    }  
    private Node fetchNode(Node curr, String prefix, int currIndex){  
        if(curr == null) return null;  
        if(currIndex == prefix.length()) return curr;  
        return fetchNode(curr.next[prefix.charAt(currIndex)], prefix, currIndex + 1);  
    }  
}
```



Tail recursion!

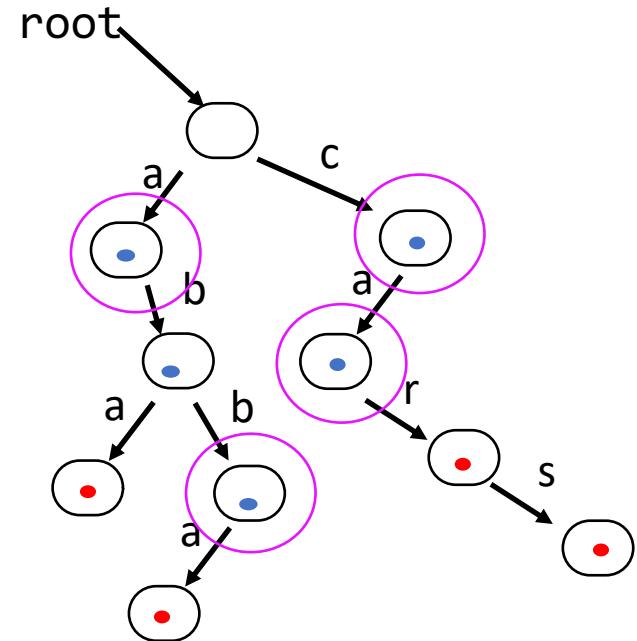
Trie redundancy

- Not everything is perfect in the trie universe!



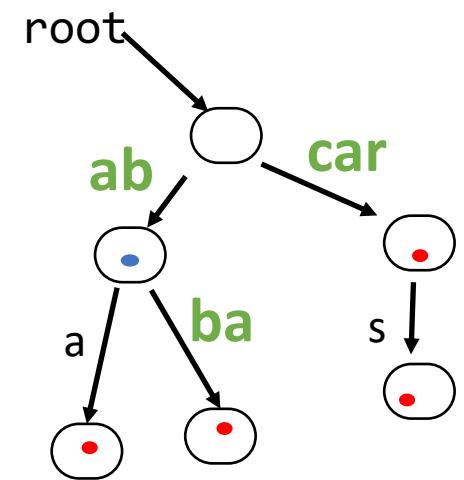
Trie redundancy

- Not everything is perfect in the trie universe!
- Do you agree with me that all the nodes that I circled on the right are effectively useless overhead?



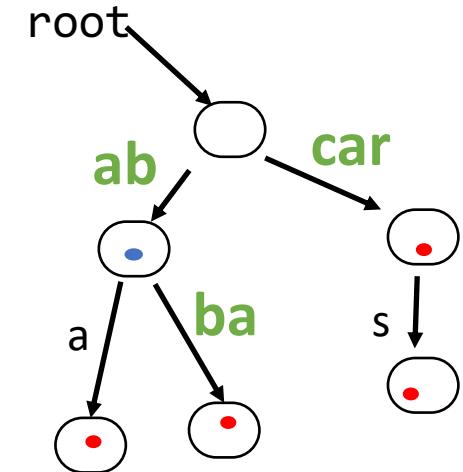
Trie redundancy

- Not everything is perfect in the trie universe!
- Do you agree with me that all the nodes that I circled on the right are effectively useless overhead?
- It would be much better if we could have something like this!



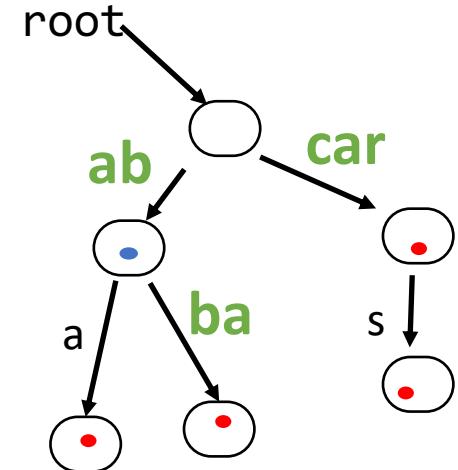
Trie redundancy

- Not everything is perfect in the trie universe!
- Do you agree with me that all the nodes that I circled on the right are effectively useless overhead?
- It would be much better if we could have something like this!
- This we will accomplish with ***Patricia Tries***, a compressed form of trie!



Trie redundancy

- Not everything is perfect in the trie universe!
- Do you agree with me that all the nodes that I circled on the right are effectively useless overhead?
- It would be much better if we could have something like this!
- This we will accomplish with ***Patricia Tries***, a compressed form of trie!



Coming Up

Next . . .