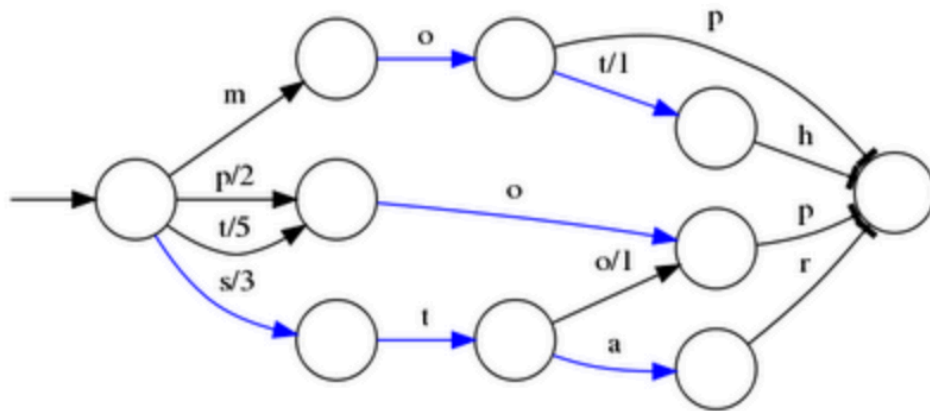


Friday, December 3, 2010

Using Finite State Transducers in Lucene

FSTs are finite-state machines that map a term (byte sequence) to an arbitrary output. They also look cool:



That FST maps the sorted words *mop*, *moth*, *pop*, *star*, *stop* and *top* to their ordinal number (0, 1, 2, ...). As you traverse the arcs, you sum up the outputs, so *stop* hits 3 on the **s** and 1 on the **o**, so its output ordinal is 4. The outputs can be arbitrary numbers or byte sequences, or combinations, etc. -- it's pluggable.

Essentially, an FST is a `SortedMap<ByteSequence,SomeOutput>`, if the arcs are in sorted order. With the right representation, it requires far less RAM than other SortedMap implementations, but has a higher CPU cost during lookup. The low memory footprint is vital for Lucene since an index can easily have many millions (sometimes, billions!) of unique terms.

Subscribe To

Posts

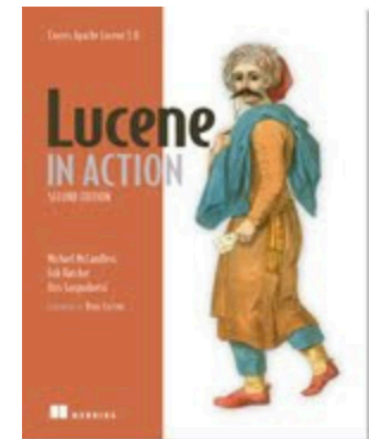
Comments

About Me

 **Michael McCandless**

[View my complete profile](#)

Lucene in Action



Blog Archive

► [2019](#) (1)

► [2017](#) (3)

► [2016](#) (1)

January, 1

Jan, 1

February, 2

Feb, 2

March, 3

Mar, 3

April, 4

Apr, 4

May, 5

June, 6

Jun, 6

July, 7

Jul, 7

August, 8

Aug, 8

September, 8

Sep, 8

October, 10

Oct, 10

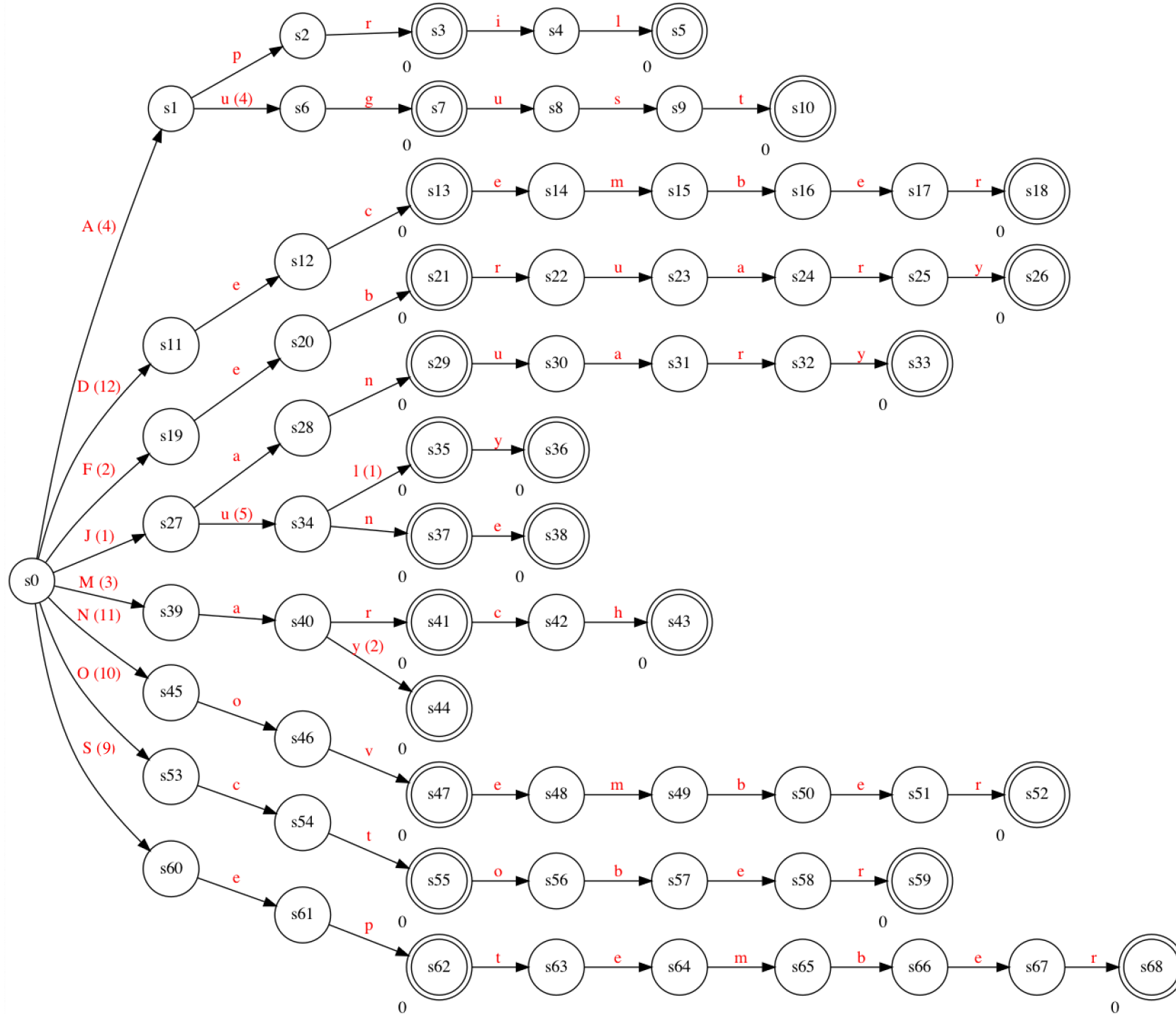
November, 11

Nov, 11

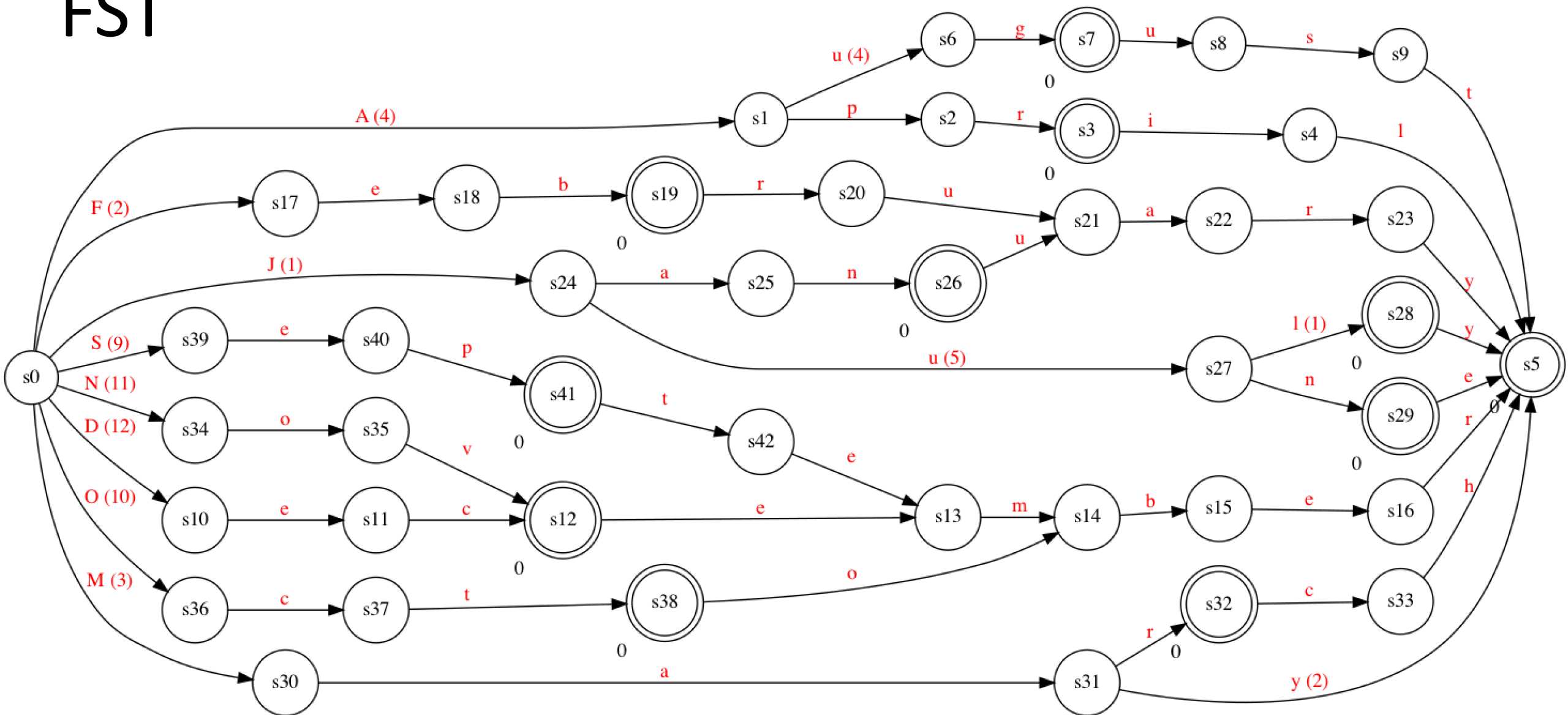
December, 12

Dec,12

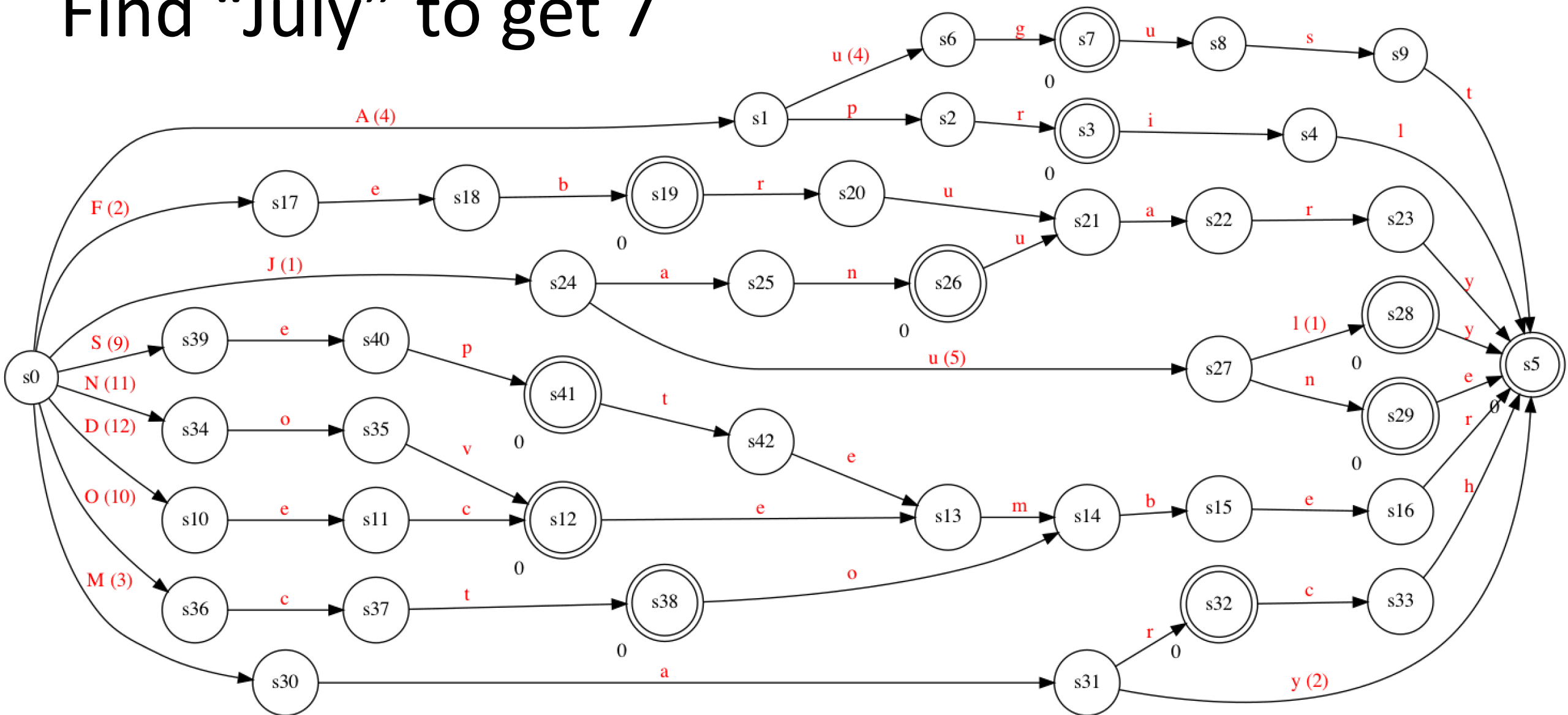
Trie



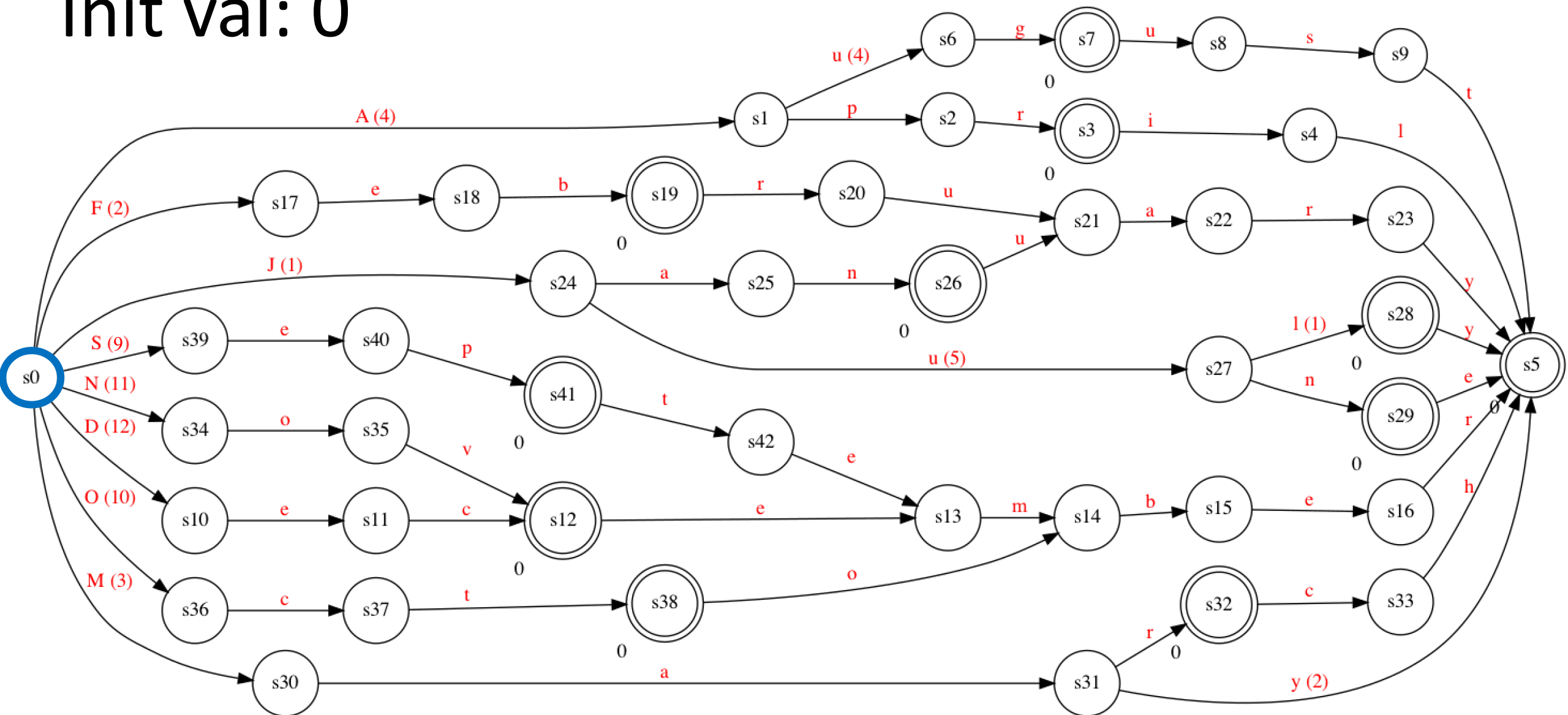
FST



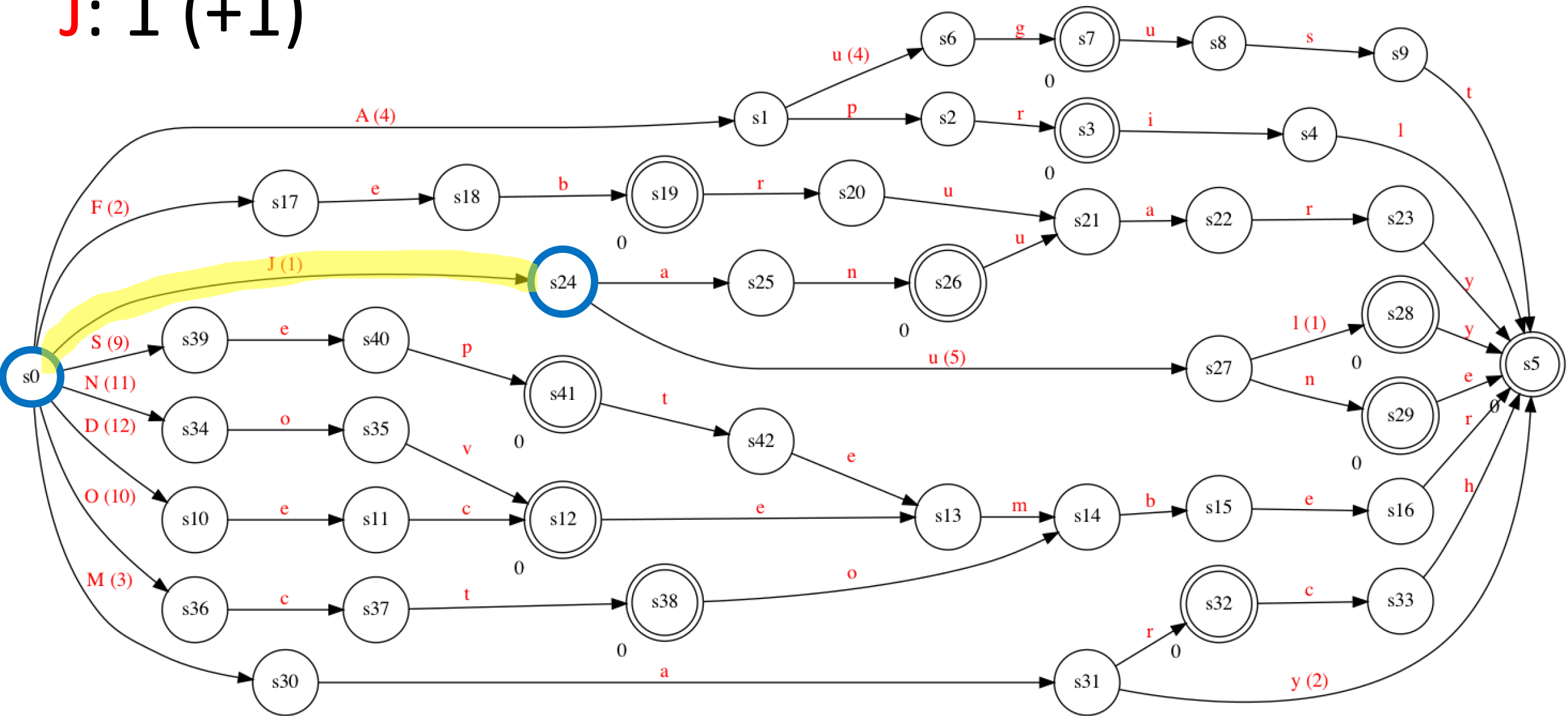
Find “July” to get 7



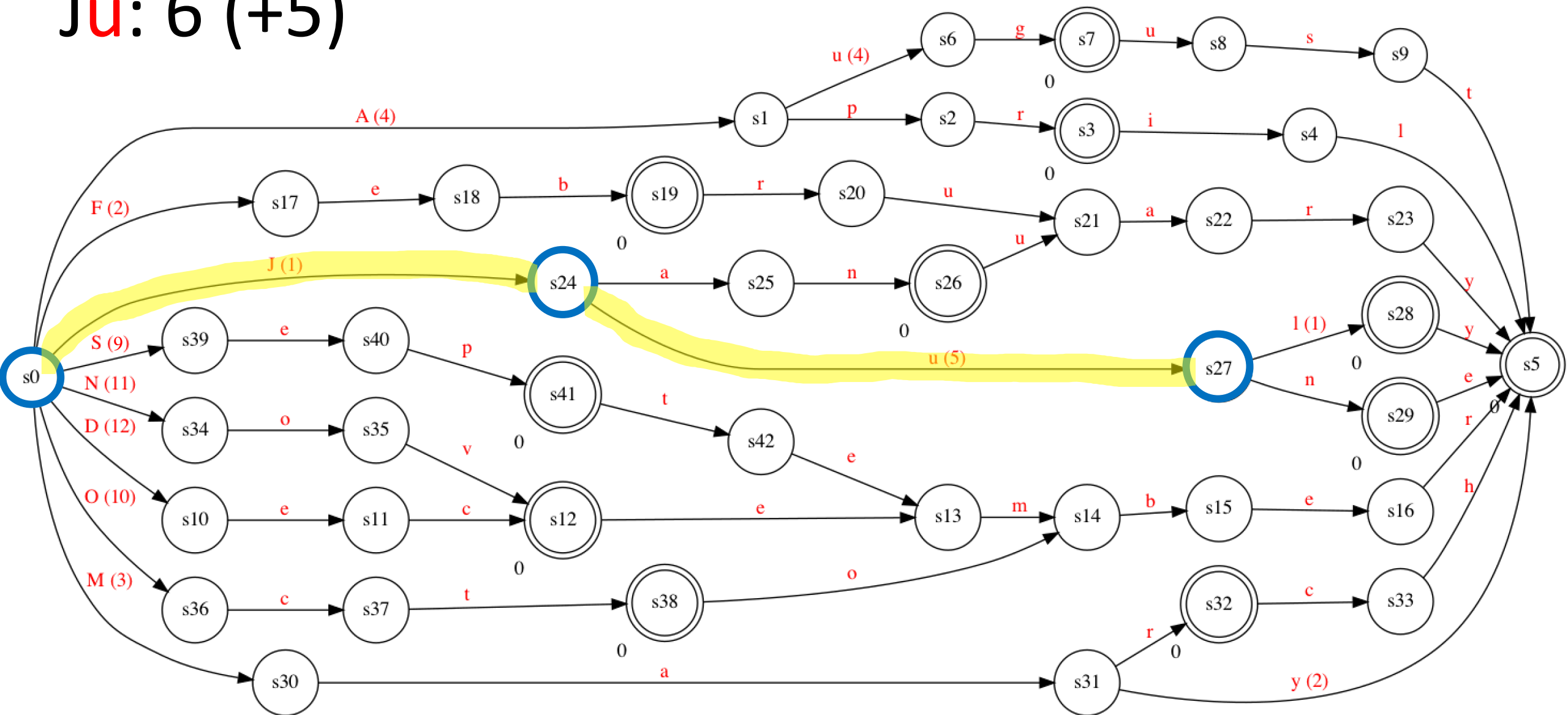
Init val: 0



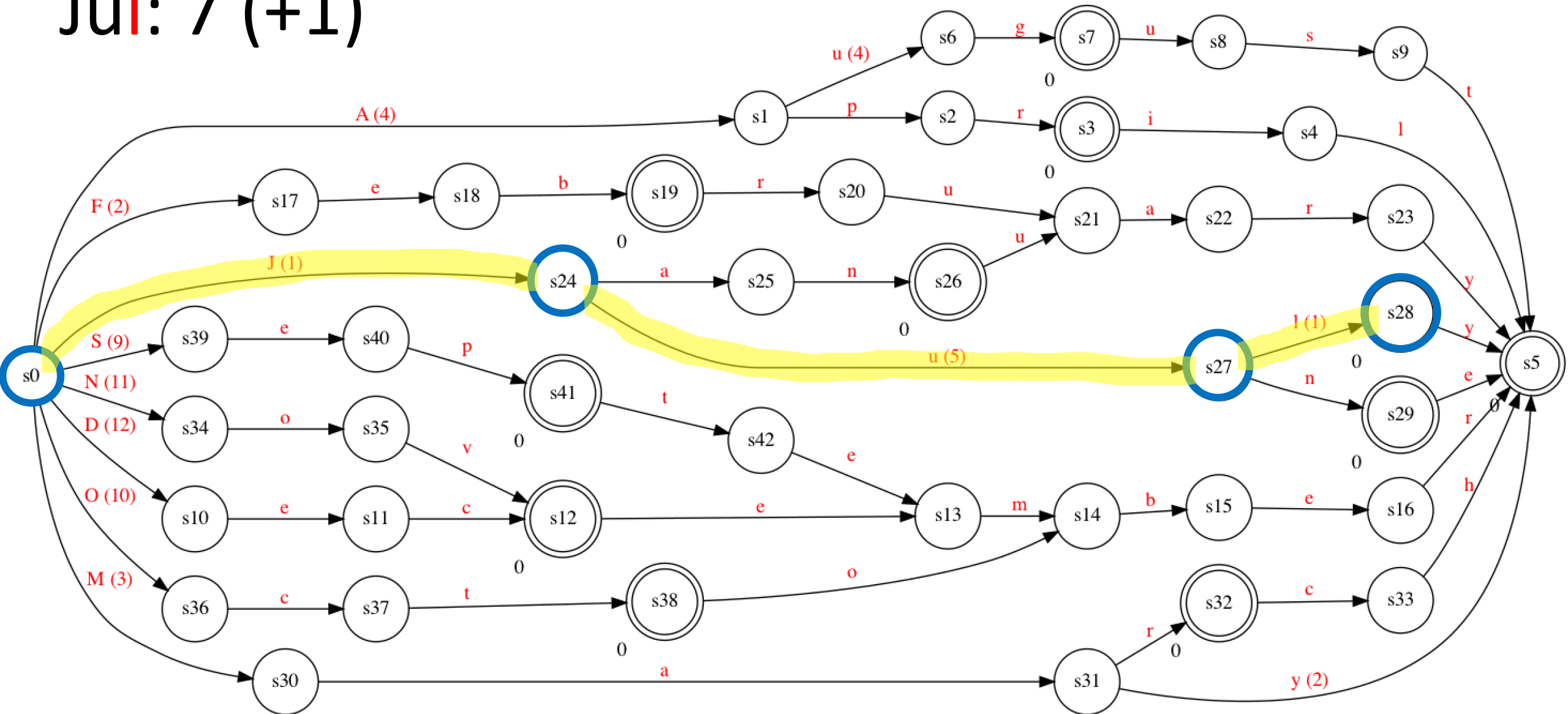
J: 1 (+1)



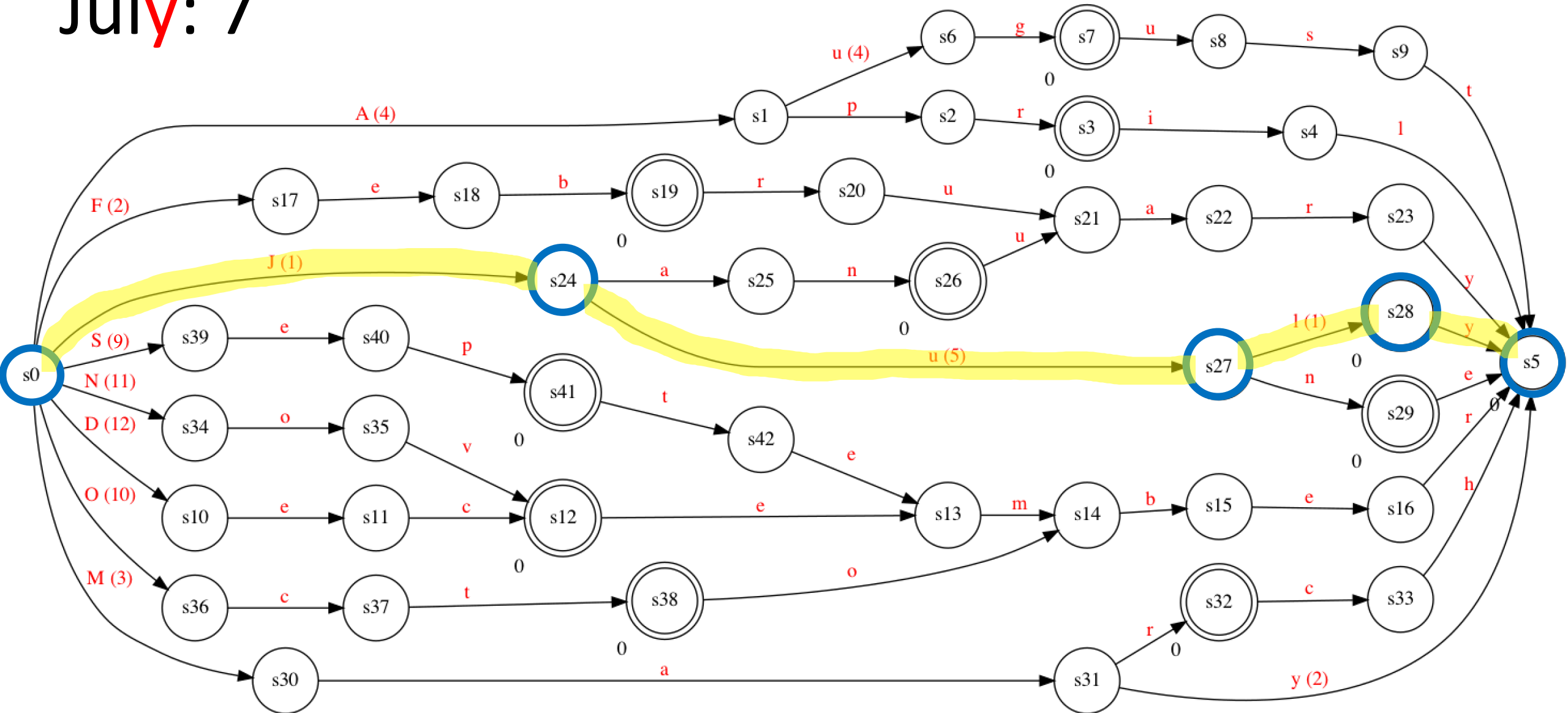
Ju: 6 (+5)



Jul: 7 (+1)



July: 7



Tree Structure

Contains pointers and containers

Can't serialize it easily to a file...

```
// Tree Structure
struct State {
    struct Transition {
        char arc;
        State* state;
        uint32_t output;
    }

    bool final;
    std::vector<Transition> transitions;
}

State* root;

// Byte Code
struct Ope {
    char arc;
    bool last_transition;
    bool final;
    uint32_t next_address;
    uint32_t output;
};

std::vector<Ope> byte_code;
```

Byte Code

Just primitive values

Can directly save it to file
Zero loading time!

Byte code

```
[master] ~/Projects/cpp-fstlib/demo$ xxd months.fst
00000000: 806c 0669 0347 8074 0673 0363 6707 846b  .l.i.G.t.s.cg..k
00000010: 8b70 0080 4623 c36d 0323 e723 8079 0643  .p..F#.m.##.y.C
00000020: 8363 43c7 2385 626e 0780 7906 8026 6e07  .cC.#.bn..y.&n.
00000030: 8185 6c0c 8c82 8583 6880 6806 e382 8079  ..l.....h.h...y
00000040: 0e85 4483 ab76 066f 03b3 6f02 7407 e3b7  ..D..v.o..o.t...
00000050: 2274 0370 0723 8953 0b8a 8b4f 088b 954e  "t.p.#.S...O...N
00000060: 0883 9e4d 0881 ada8 82c4 4608 8cd1 4408  ...M.....F...D.
00000070: 84de 4108 0004 080c 0f13 171b 88fe 0000  ..A.....
00000080: 0100 7d00 0000 0065 7275 614a 6263 0000  ..}....eruaJbc..
00000090: 0000 0000 0000  ....
```

Byte code dump

Byte code records are stored in backward order!

- N: No next address. Just move to next record
- F: Final state (or ‘accepted’)
- L: Last transition in a state

Address	Arc	N	F	L	NxtAddr	Output	StOuts	Size
2	l		*	-	x			3
4	i	↑	*	-				2
5	r	↑	*	-				1
8	t		*	-	x			3
10	s	↑		-				2
11	u	↑		-				1
13	g	↑	*	-				2
15	u	↑		-		4		2
18	p			-	5			3
20	r		*	-	x			2
21	e	↑		-				1
22	b	↑		-				1
24	m	↑		-				2
25	e	↑		-				1
26	c	↑	*	-				1
27	e	↑		-				1
30	y		*	-	x			3
31	r	↑		-				1
32	a	↑		-				1
33	u	↑		-				1
34	r	↑		-				1
35	b	↑	*	-				1
36	e	↑		-				1
38	u			-	32			2
40	n	↑	*	-				2
43	y		*	-	x			3
45	e		*	-	x			2
47	n	↑	*	-				2
51	a		*	-	43	1		4
53	a			-	40			2
56	u			-	51	5		3
59	h		*	-	x			3
60	c	↑		-				1
64	y		*	-	x	2		4
66	r		*	-	60			2
67	a	↑		-				1
70	v		*	-	25			3
72	o	↑		-				2
75	o			-	22			3
77	t	↑	*	-				2
78	c	↑		-				1
80	e			-	24			2
82	t	↑		-				2
84	p	↑	*	-				2
85	e	↑		-				1
88	S	↑		-		9		3
92	O			-	78	10		4
96	N			-	72	11		4
100	M			-	67	3		4
103	J			-	56	1		3
107	F			-	36	2		4
111	D			-	27	12		4
125	A			-	18	4		14

FST (Finite State Transducer) library

Trie-like data structure

$O(M)$ access like `std::unordered_map`

Common prefix match ('computer programming is fun.' returns 'computer', 'computer programming')

Read-only binary database

Zero loading time when lookup

Memory-mapped file friendly

Compact ('/usr/shared/dict/words': 2.4M → 1.1M)

Examples that fstlib can do, but `std::unordered_map` cannot effectively

Keyword extraction ('Longest common prefix search')

Word/phrase list dictionary for full-text search library ('Zero loading time', 'Compact')

cpp-fstlib internals

Compile a source dictionary to byte code, then interpret it when searching

Incremental direct database construction without having a complete set of states on memory

Implemented based on <https://blog.burntsushi.net/transducers/>. Thanks for the valuable information!