# Range Filters

(SuRF, Memento Filter, Diva)

Niv Dayan - CSC2525: Research Topics in Database Management

# Last lecture



This was fun...

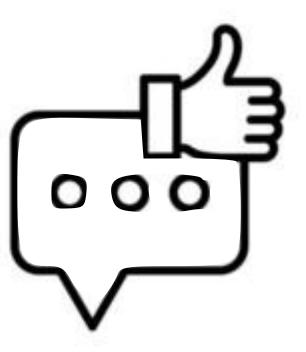
# Logistics



Projects due on April 9



Oral Exams April 8-10



Course Evaluation

# **Course Evaluation**

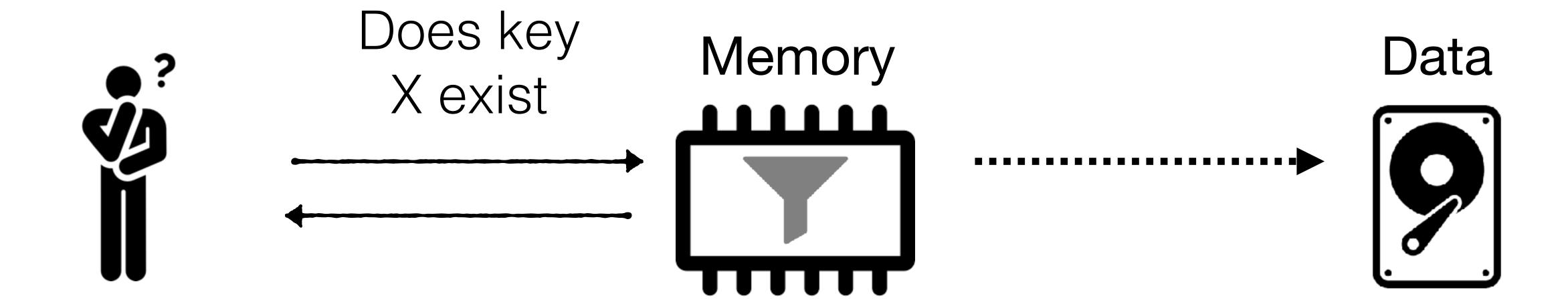


Helps us improve

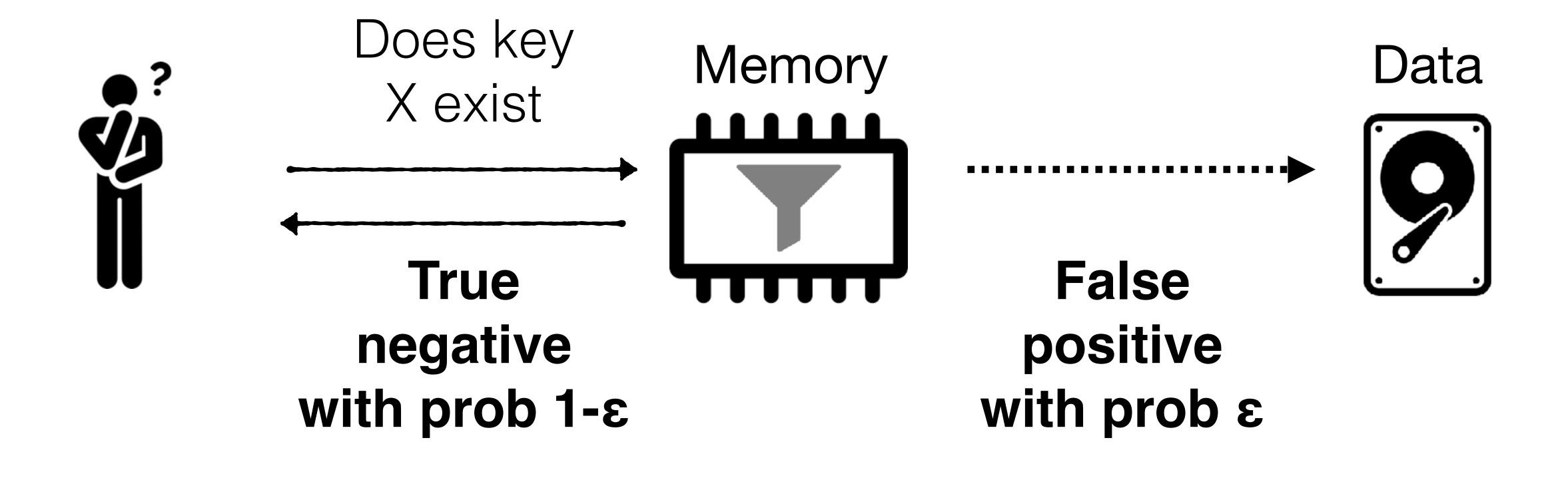
# What is a Filter?

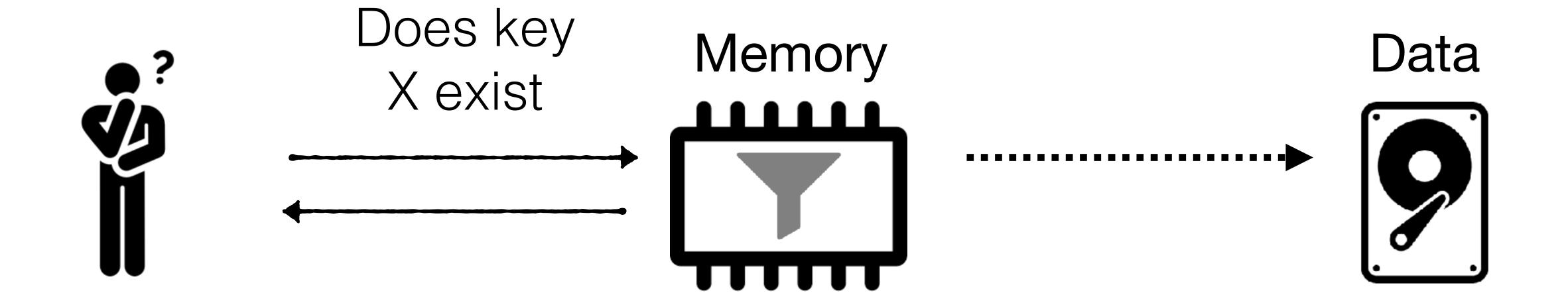
Does X exist? → X Y Z

# If key X does not exist



# If key X does not exist





Saves storage accesses & network hops

# only support point queries (to one key)



# only support point queries (to one key)

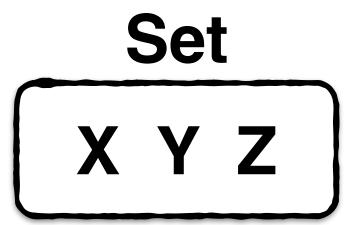
# How about a range?



How about a range?

Does anything in [A, B] exist?

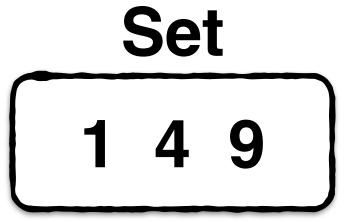




How about a range?

Does anything in [3, 5] exist?





# How about a range?

Does anything in [3, 5] exist?

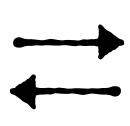
True positive







Set
1 4 9





True negative with prob 1-ε

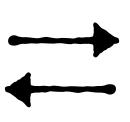


True negative with prob 1-ε

False positive with prob ε

#### **Applications?**

Does anything in [5, 8] exist?

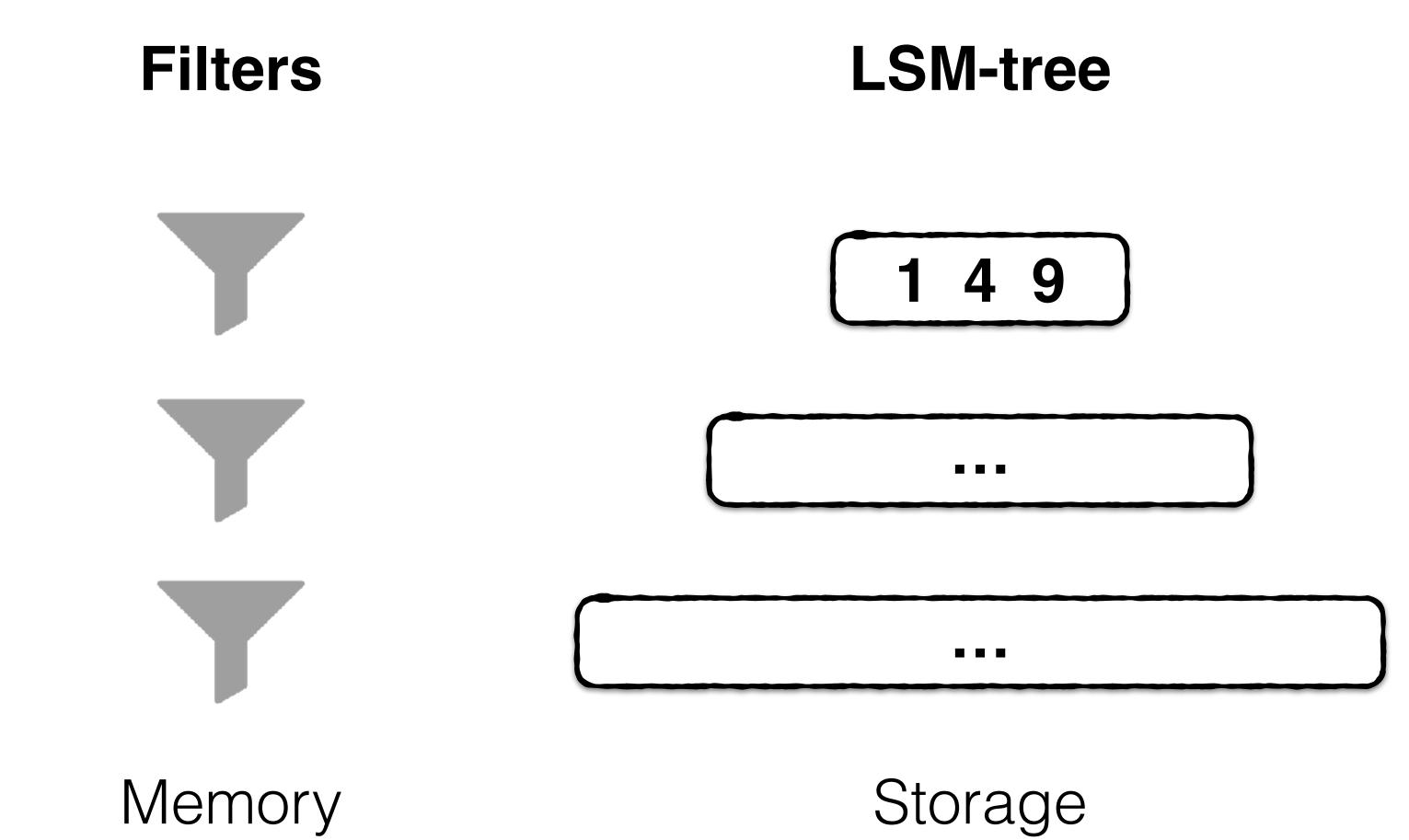


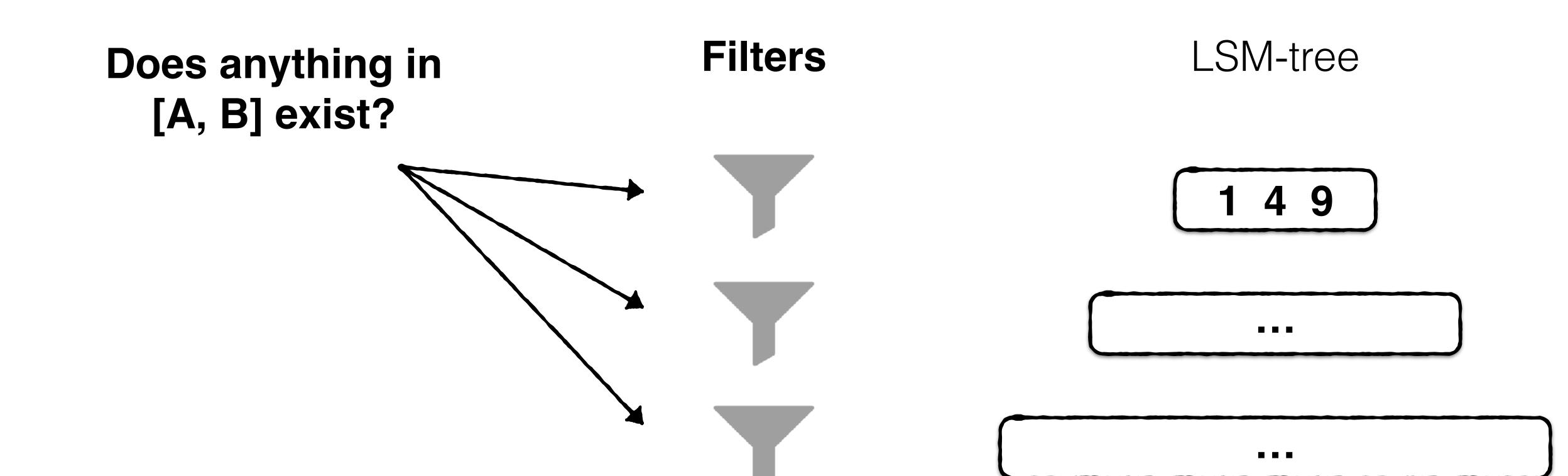




True negative with prob 1-ε

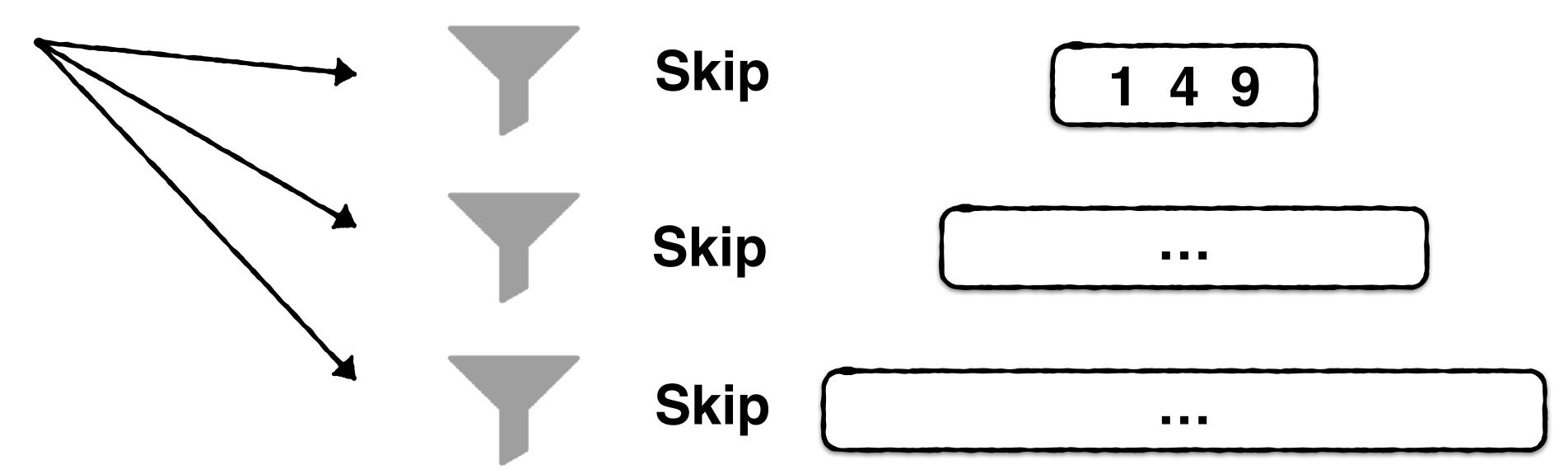
False positive with prob ε





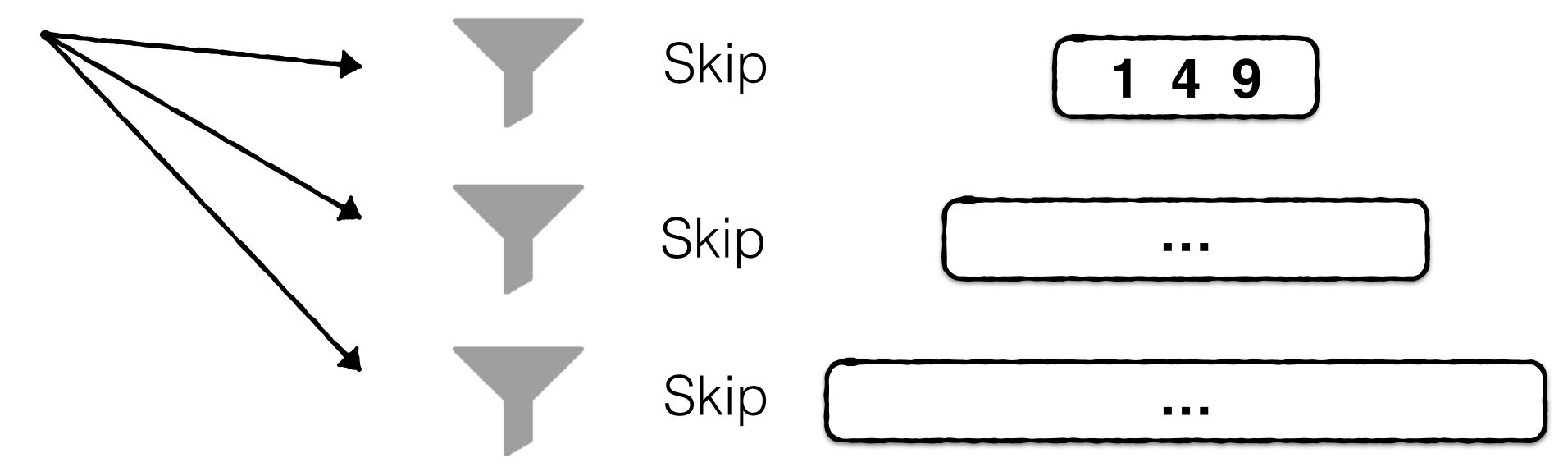
## Ideally skip runs that don't contain relevant key range

Does anything in [A, B] exist?



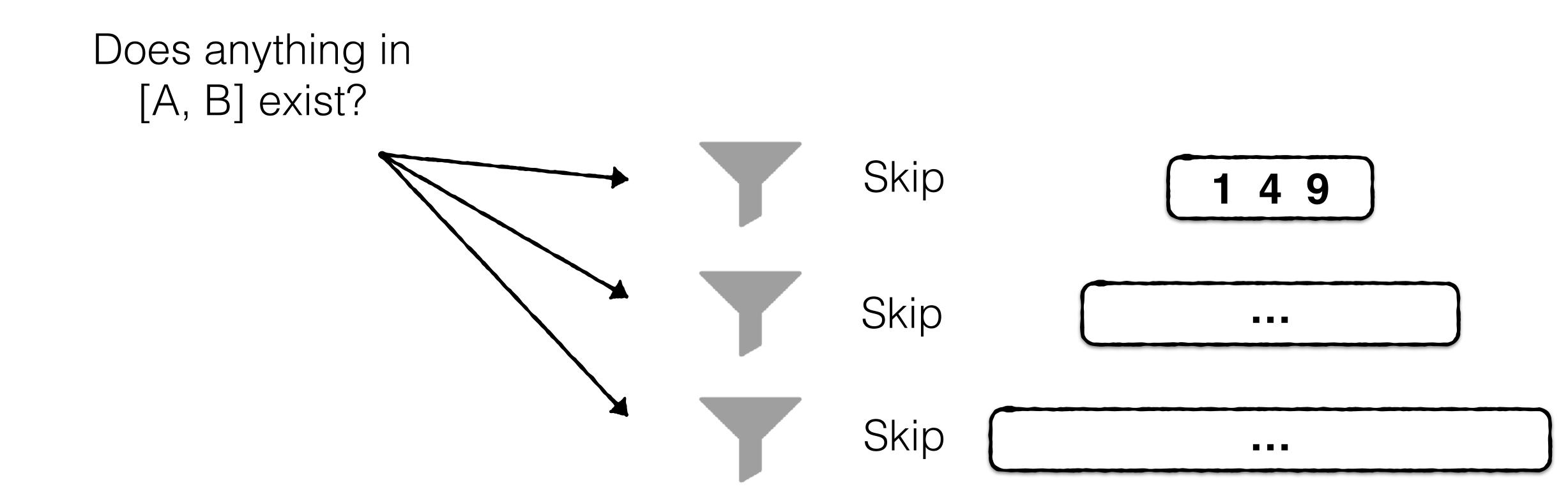
# How to implement a range filter?

Does anything in [A, B] exist?



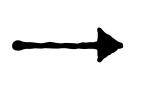
#### How to implement a range filter?

## Can we use a point filter (Bloom, Quotient) to answer range queries?



#### **Bloom**

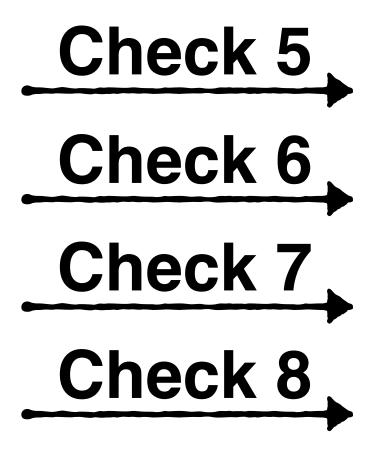
Does anything in [5, 8] exist?

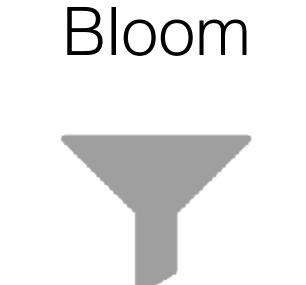




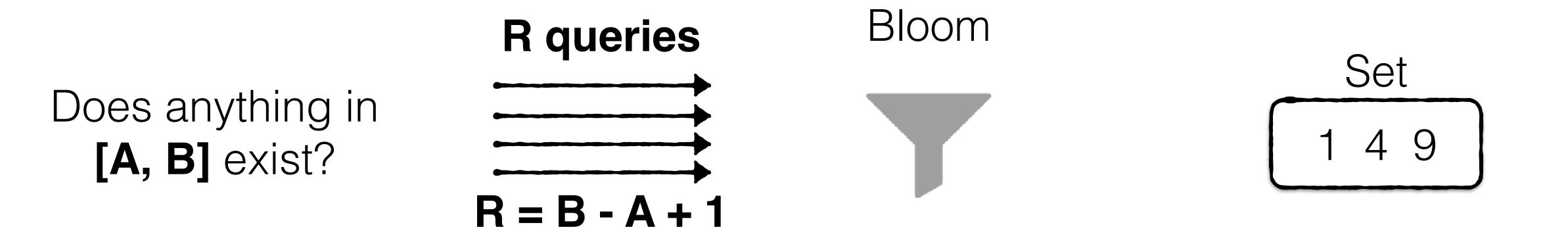


Does anything in [5, 8] exist?

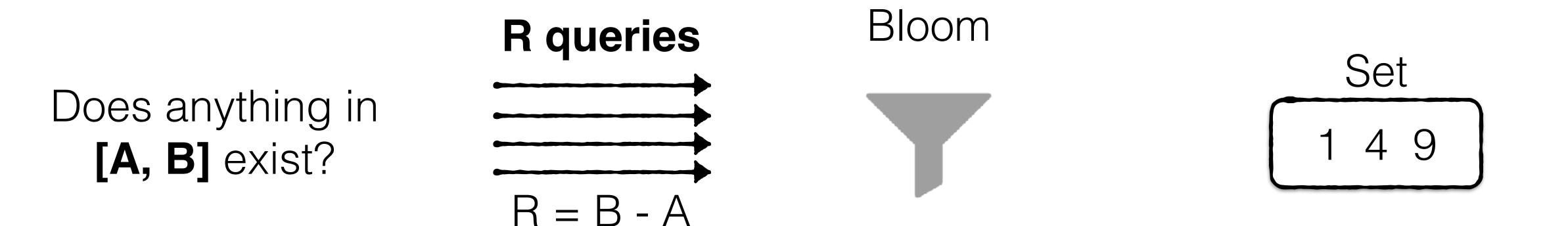








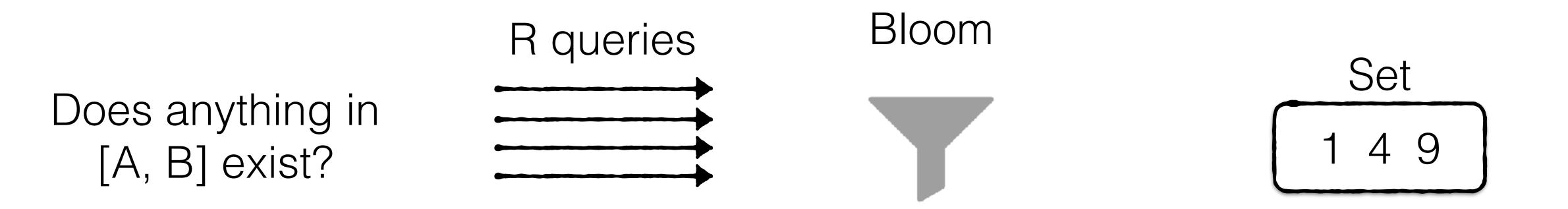
If all queries return negative, we know key does not exist



If all queries return negative, we know key does not exist

If at least one returns a positive, the overall outcome is a positive

#### **Problems?**

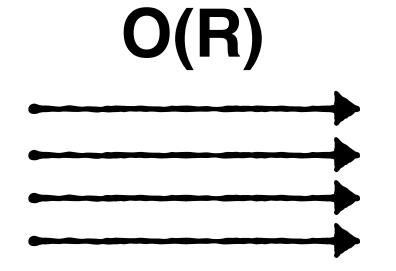


If all queries return negative, we know key does not exist

If at least one returns a positive, the overall outcome is a positive

# **Problem 1: Query Cost**

Does anything in [A, B] exist?

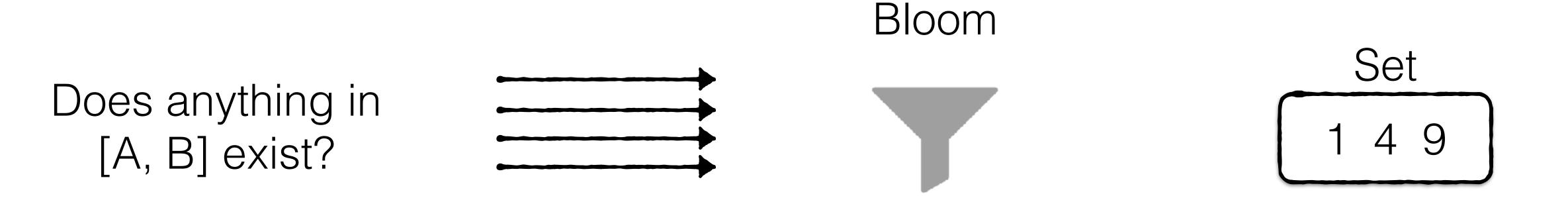


Bloom



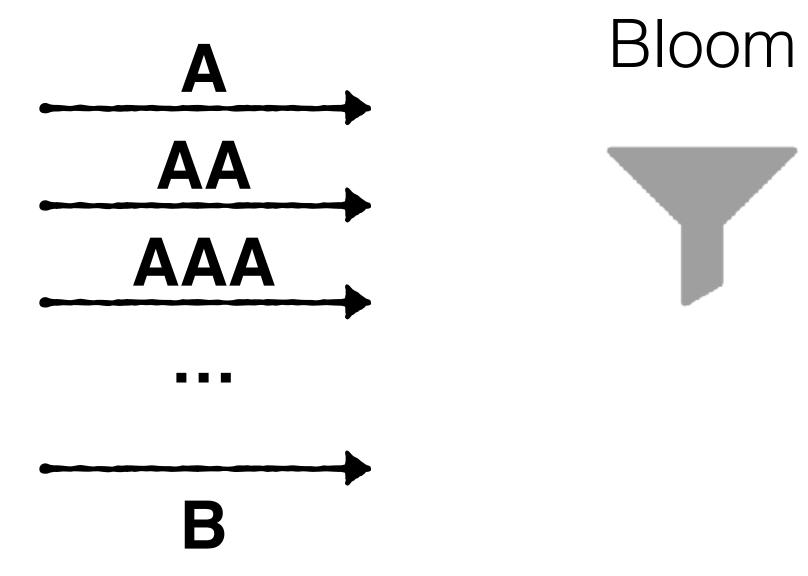
Set 1 4 9

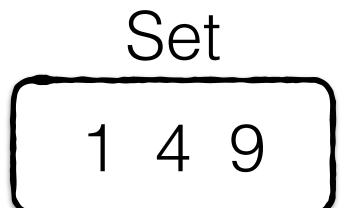
# Problem 2: Does not work for infinite universe (e.g., strings)



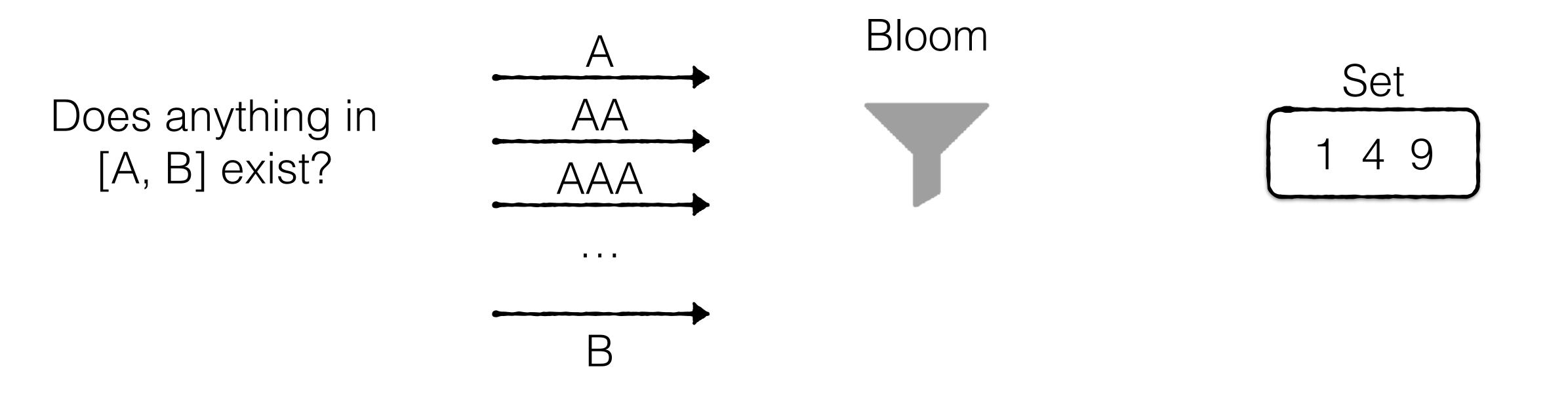
## Problem 2: Does not work for infinite universe (e.g., strings)

Does anything in [A, B] exist?



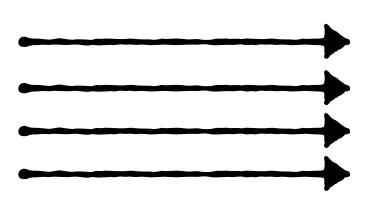


Problem 2: Does not work for infinite universe (e.g., strings)



There are infinitely many possible keys in a range:)

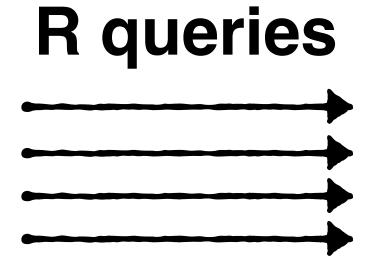
Does anything in [A, B] exist?



Bloom



Does anything in [A, B] exist?







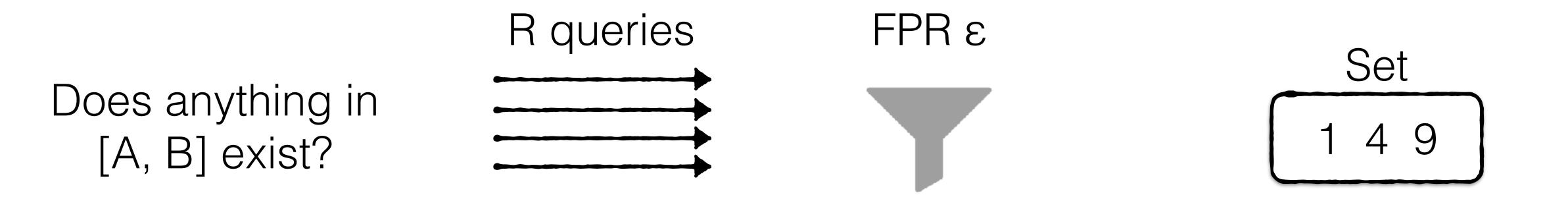
Does anything in [A, B] exist?

FPR ε

Set

1 4 9

P[false positive]?



P[false positive] = P[at least one query returns positive]

Does anything in [A, B] exist?

R queries FPR ε

Set

1 4 9

P[false positive] = P[at least one query returns positive]

= 1 - P[all queries return negative]

Does anything in [A, B] exist?

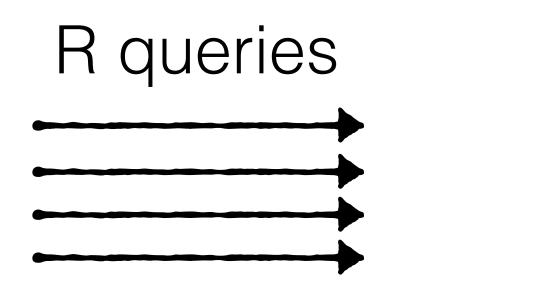


 $P[false\ positive] = P[at\ least\ one\ query\ returns\ positive]$ 

= 1 - P[all queries return negative]

= 1 - P[one query return negative]<sup>R</sup>

Does anything in [A, B] exist?







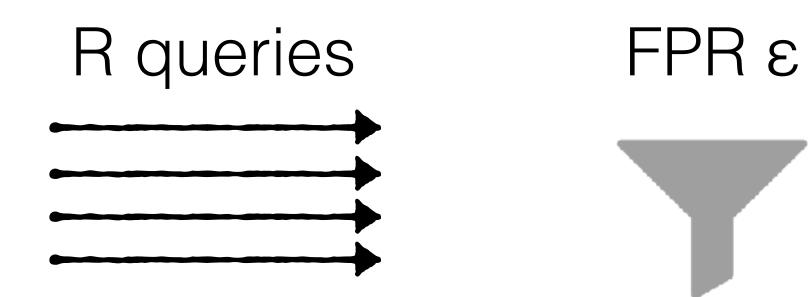


```
P[false positive] = P[at least one query returns positive]
```

- = 1 P[all queries return negative]
- = 1 P[one query return negative]R

$$= 1 - (1-\varepsilon)^{R}$$

Does anything in [A, B] exist?



```
Set
1 4 9
```

```
P[false positive] = P[at least one query returns positive]

= 1 - P[all queries return negative]

= 1 - P[one query return negative]<sup>R</sup>

= 1 - (1-ε)<sup>R</sup>

≈ 1 - eε·R
```

Does anything in [A, B] exist?



P[false positive] = P[at least one query returns positive] = 1 - P[all queries return negative] = 1 - P[one query return negative]<sup>R</sup> = 1 -  $(1-\epsilon)^R$   $\approx 1 - e^{\epsilon \cdot R}$ <  $\epsilon \cdot R$  By union bound

Set

4 9

Does anything in [A, B] exist?

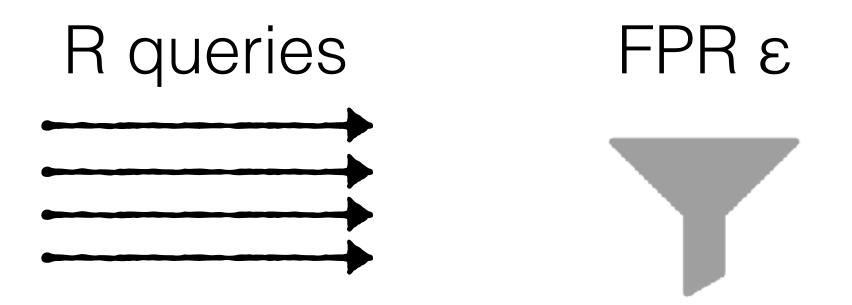
R queries FPR ε

Set

1 4 9

P[false positive]  $< \epsilon \cdot R$ 

Does anything in [A, B] exist?



Set

1 4 9

P[false positive]  $< \epsilon \cdot R$ 

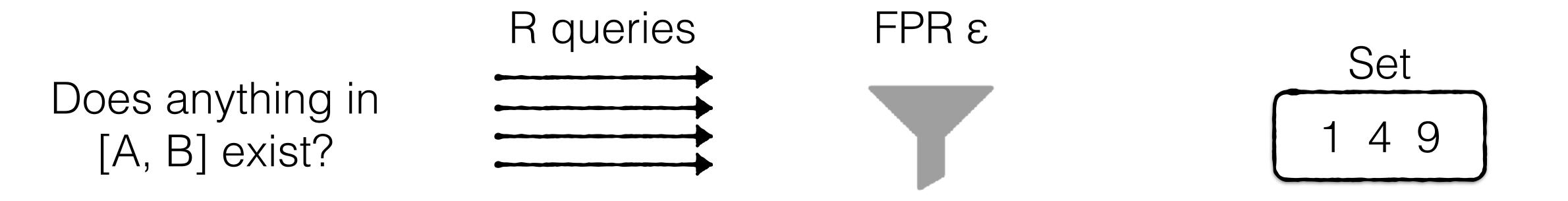
**Useless for large R** 

Does anything in [A, B] exist?



Recall that:

$$\varepsilon = 2^{-F}$$
 where  $F = bits / entry$ 



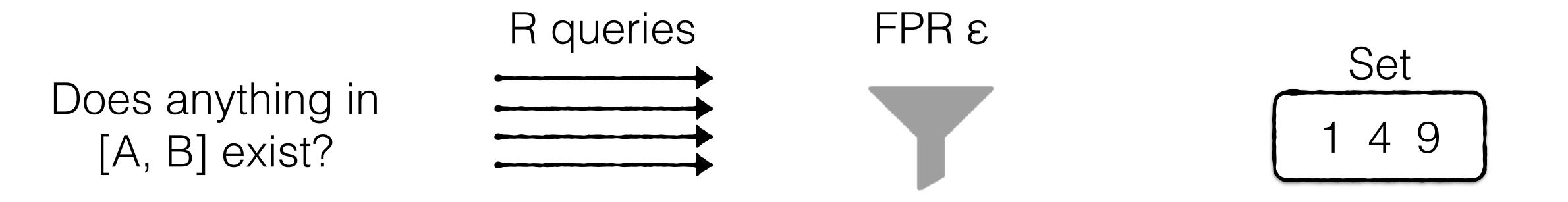
P[false positive] < 2-F · R

How many extra bits per entry do we need to make up for R?

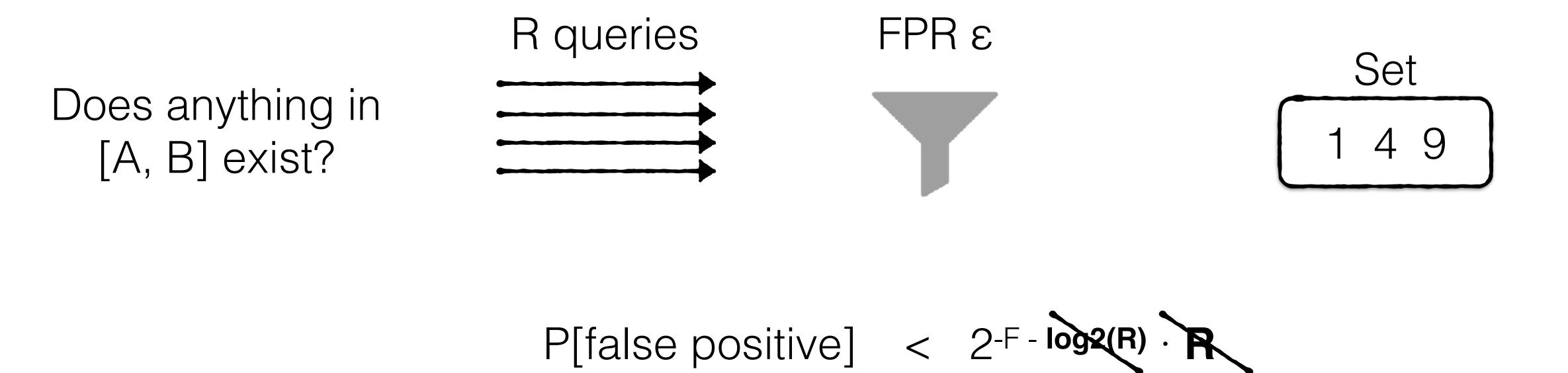


P[false positive] < 2-F · R

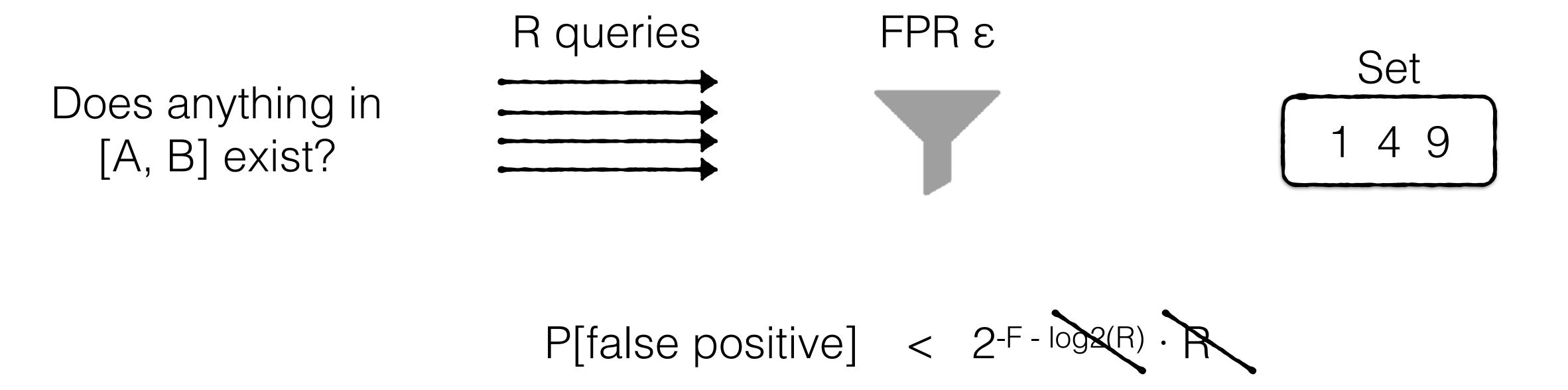
How many extra bits per entry do we need to make up for R?  $log_2(R)$ 



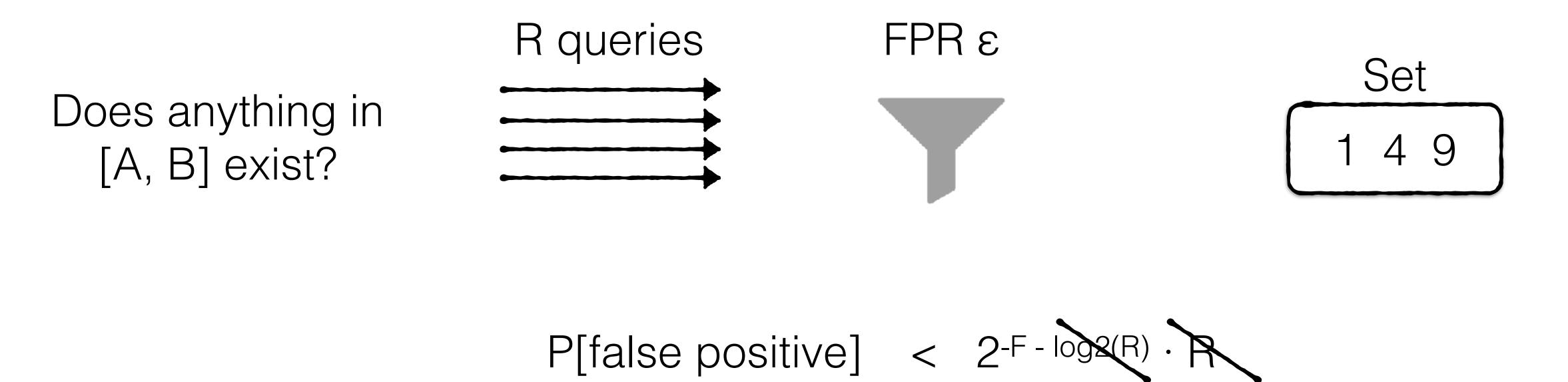
How many extra bits per entry do we need to make up for R?



How many extra bits per entry do we need to make up for R?



FPR is higher by factor of R, or we need log<sub>2</sub>(R) extra bits / entry to keep it stable



FPR is higher by factor of R, or we need log<sub>2</sub>(R) extra bits / entry to keep it stable

But this also requires us to know max range query length in advance

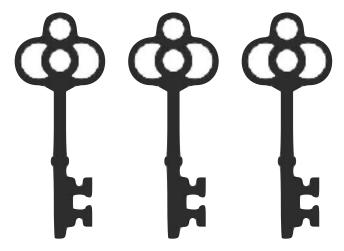
#### **Problems**



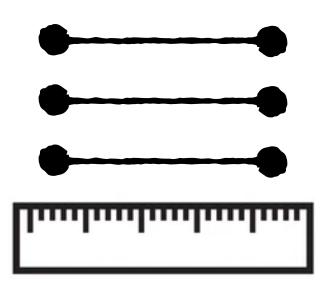
query cost O(R)



FPR O(ε·R)



Fixed-length keys



#### Problems

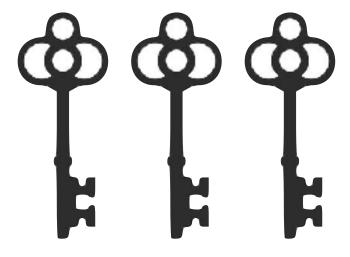
# How much can we improve these?



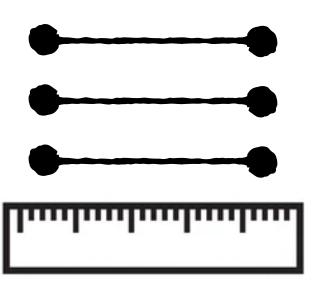
query cost O(R)



FPR O(ε·R)



Fixed-length keys





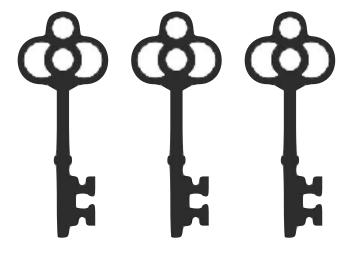
query cost O(1)



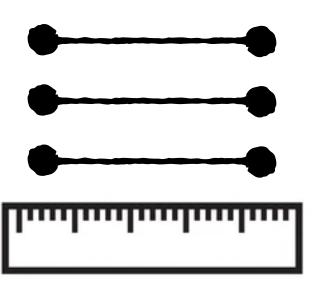
?



FPR  $O(\epsilon \cdot R)$ 



Fixed-length keys



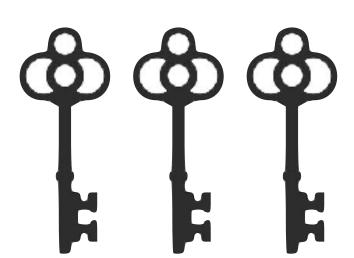


query cost O(1)

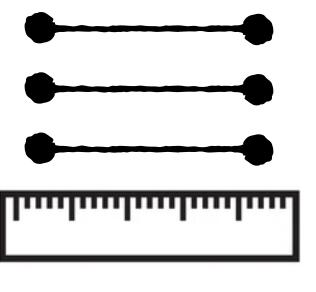


FPR O(ε·R)





Fixed-length keys

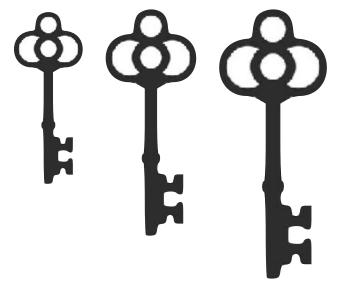




query cost O(1)

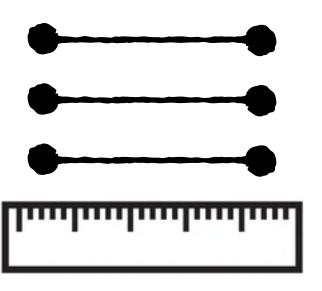


FPR O(ε)



Var-length keys



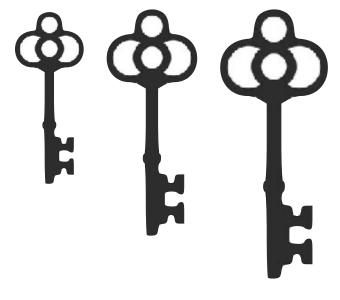




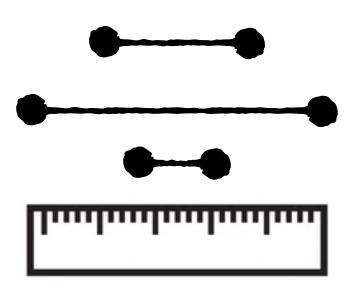
query cost O(1)



FPR O(ε)



Var-length keys



Var-length queries

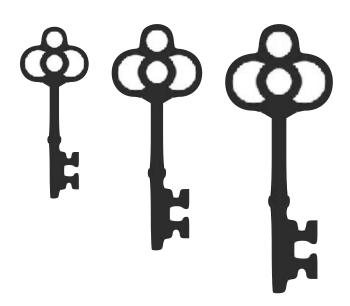




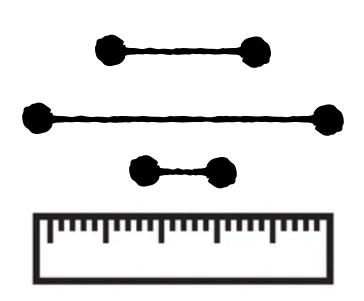
query cost O(1)



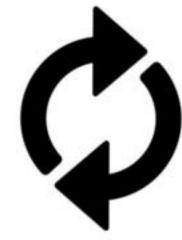
FPR O(ε)



Var-length keys



Var-length queries



+ Dynamic operations (Inserts, deletes, expansions, contractions)

#### Hot research topic in past decade

GGLRSuRFRosettaSNARFProteusSODA14SIGMOD18SIGMOD2020SIGMOD22SIGMOD22

**REncoder BloomRF Grafite Oasis Memento Diva**ICDE23 EDBT23 SIGMOD24 VLDB24 SIGMOD24 ...25

SuRF SIGMOD18

Memento SIGMOD24

Diva ...25



















Var-length keys & queries

FPR guarantee

**Query speed** 

**Dynamic** 

# SuRF

Memento

Diva







Var-length keys & queries

Yes

FPR guarantee

None

Query speed

O(L) (L = key length)

Dynamic

No

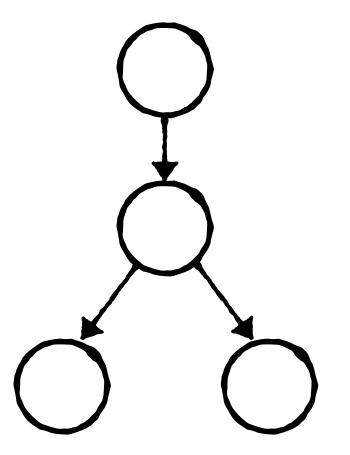
SuRF Memento Diva Var-length keys Yes No & queries None Robust FPR guarantee **O**(1) O(L)Query speed (L = key length)**Dynamic** No Yes

	SuRF	Memento	Diva
	Ž.		
Var-length keys & queries	Yes	No	Yes
FPR guarantee	None	Robust	Semi-Robust
Query speed	O(L) (L = key length)	O(1)	O(log L)
Dynamic	No	Yes	Yes

### SuRF: Practical Range Query Filtering with Fast Succinct Tries

Huanchen Zhang, Hyeontaek Lim, Viktor Leis, David G. Andersen, Michael Kaminsky, Kimberly Keeton, Andrew Pavlo

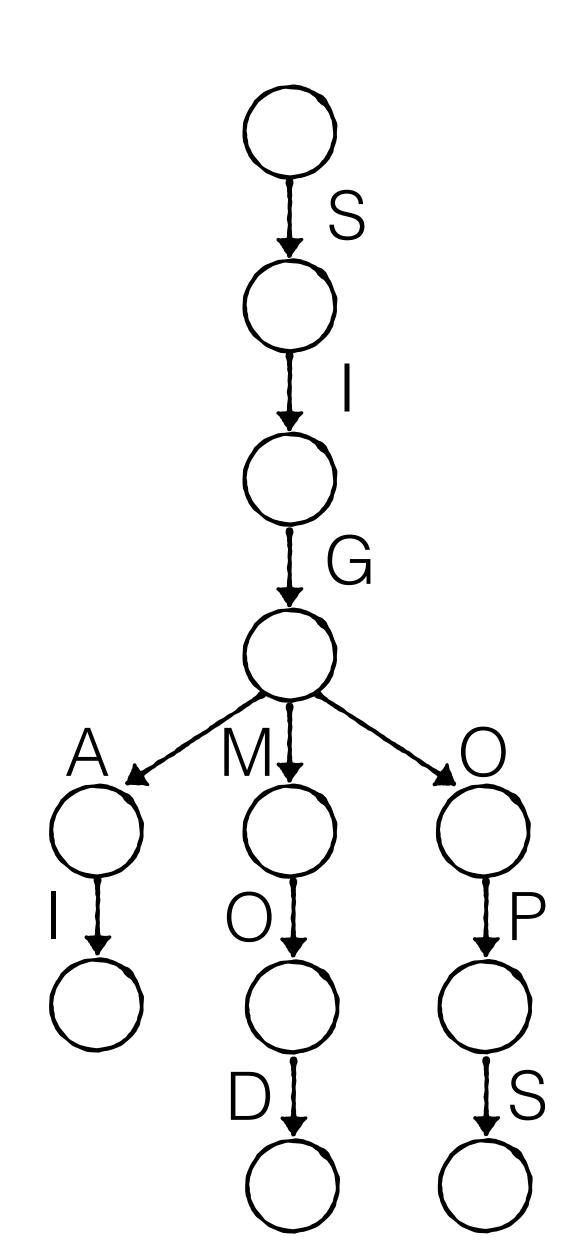
#### SIGMOD18



# **Starting Point**

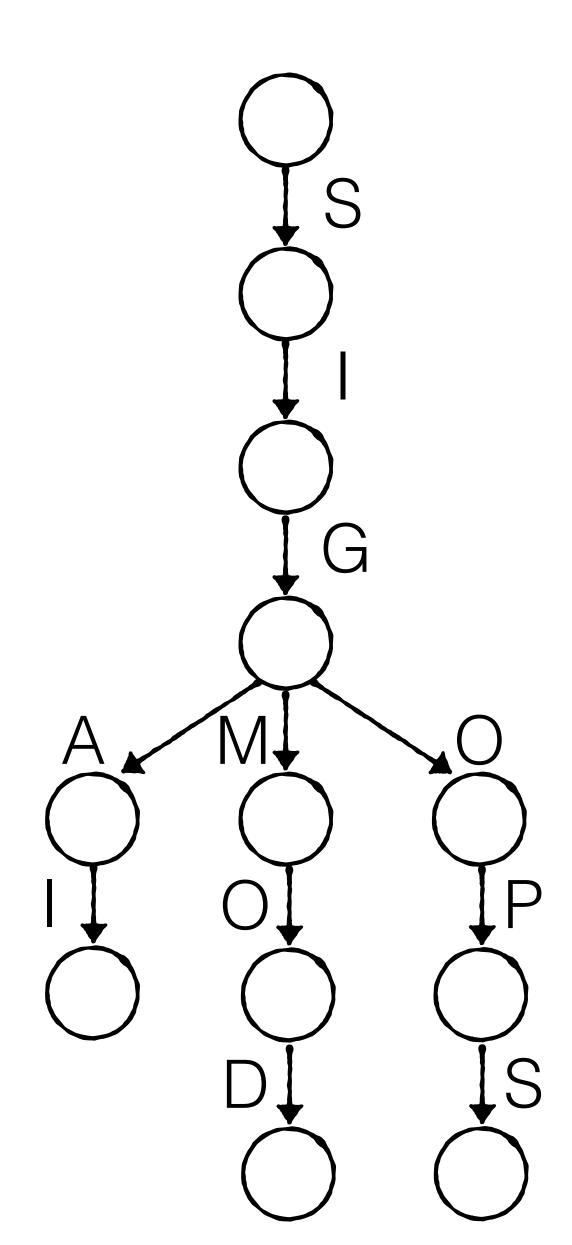
# Keys

SIGAI SIGMOD SIGOPS



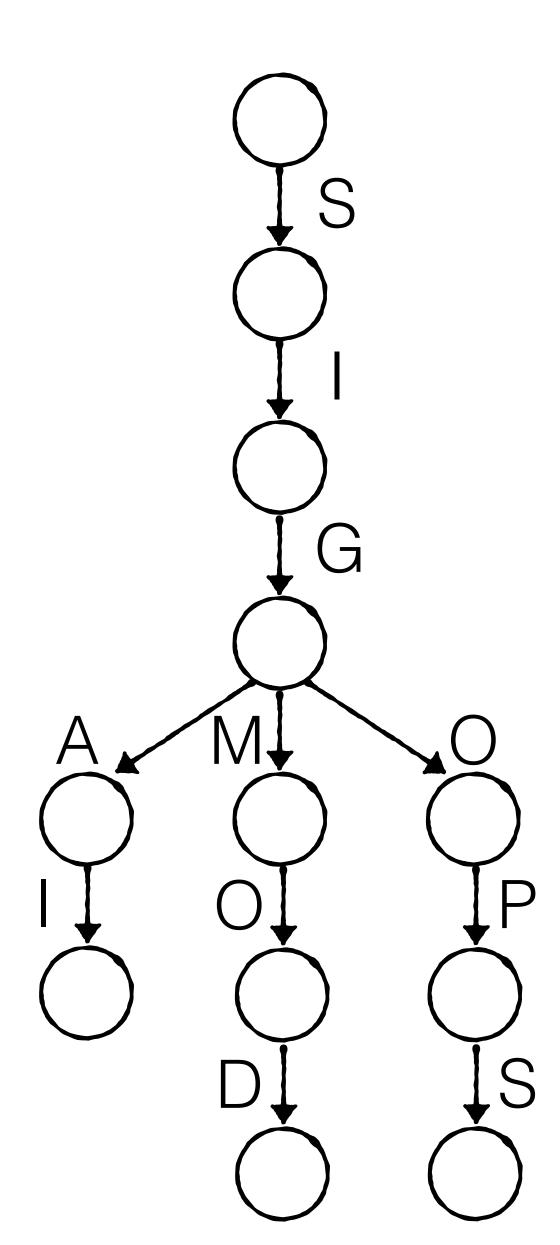
# Starting Point

Keys
SIGAI
SIGMOD
SIGOPS



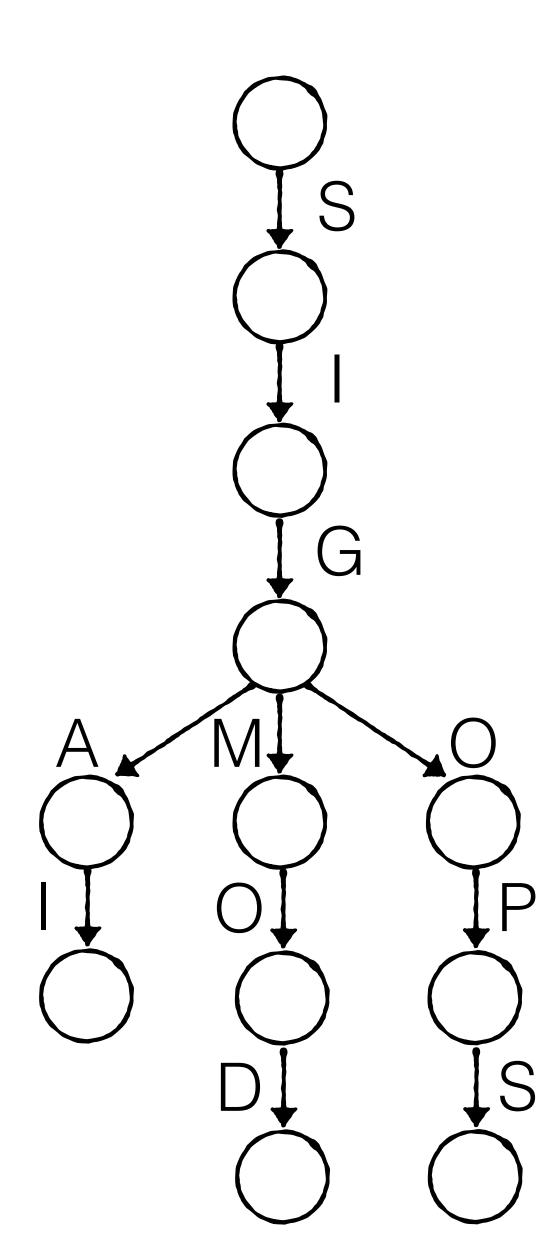
Fanout - 256 (1 byte)

# Problems?



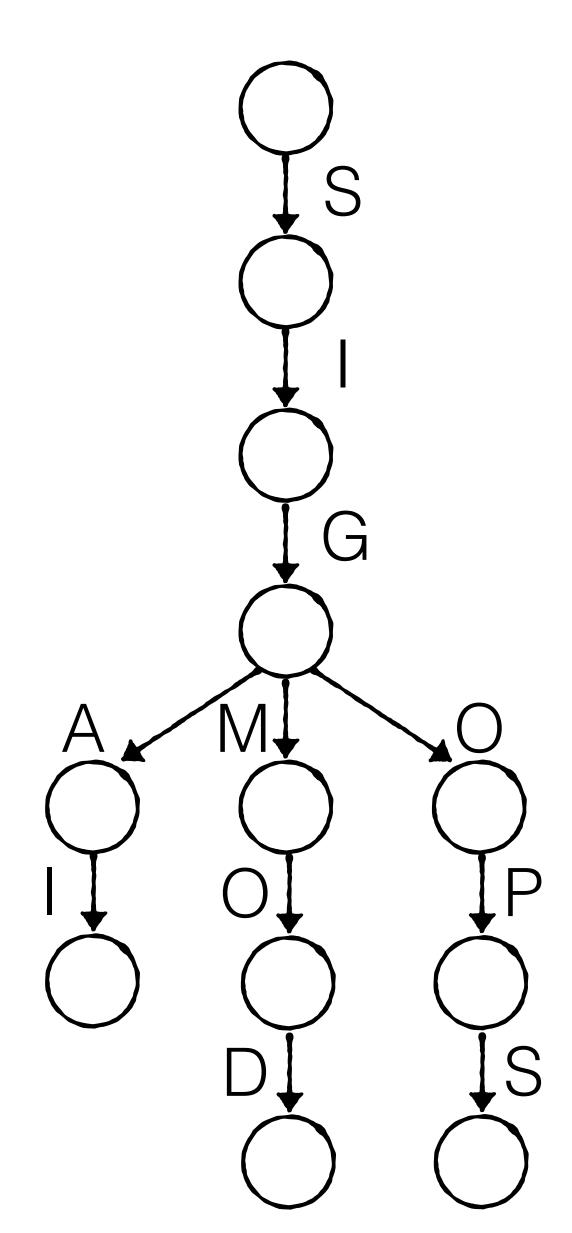
#### Problems?

# 1. Stores full keys



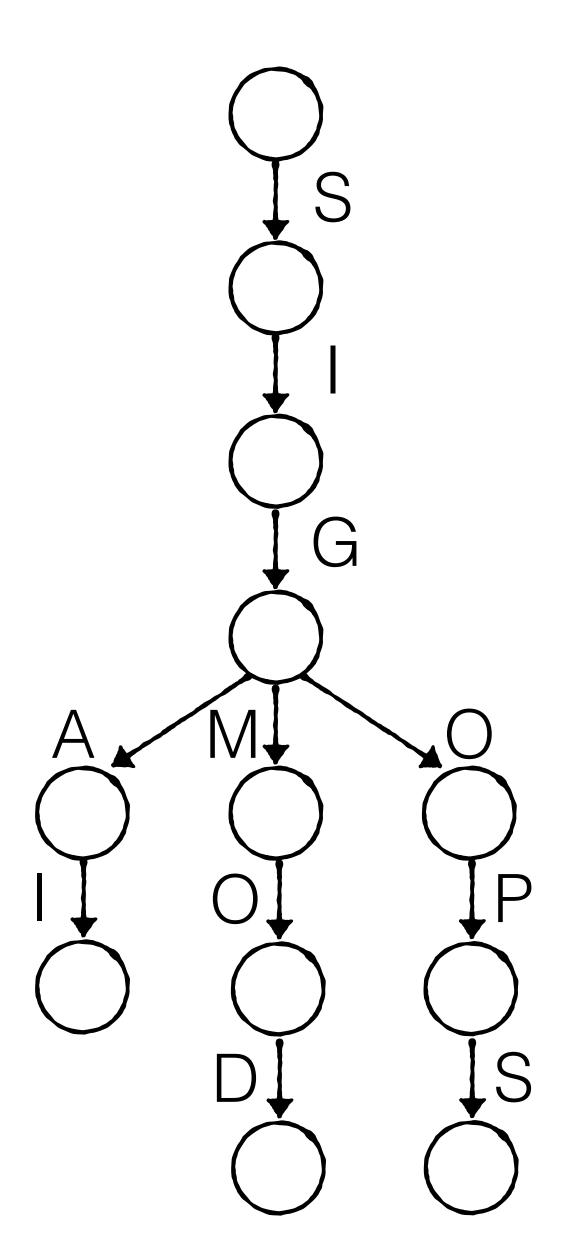
#### Problems?

1. Stores full keys

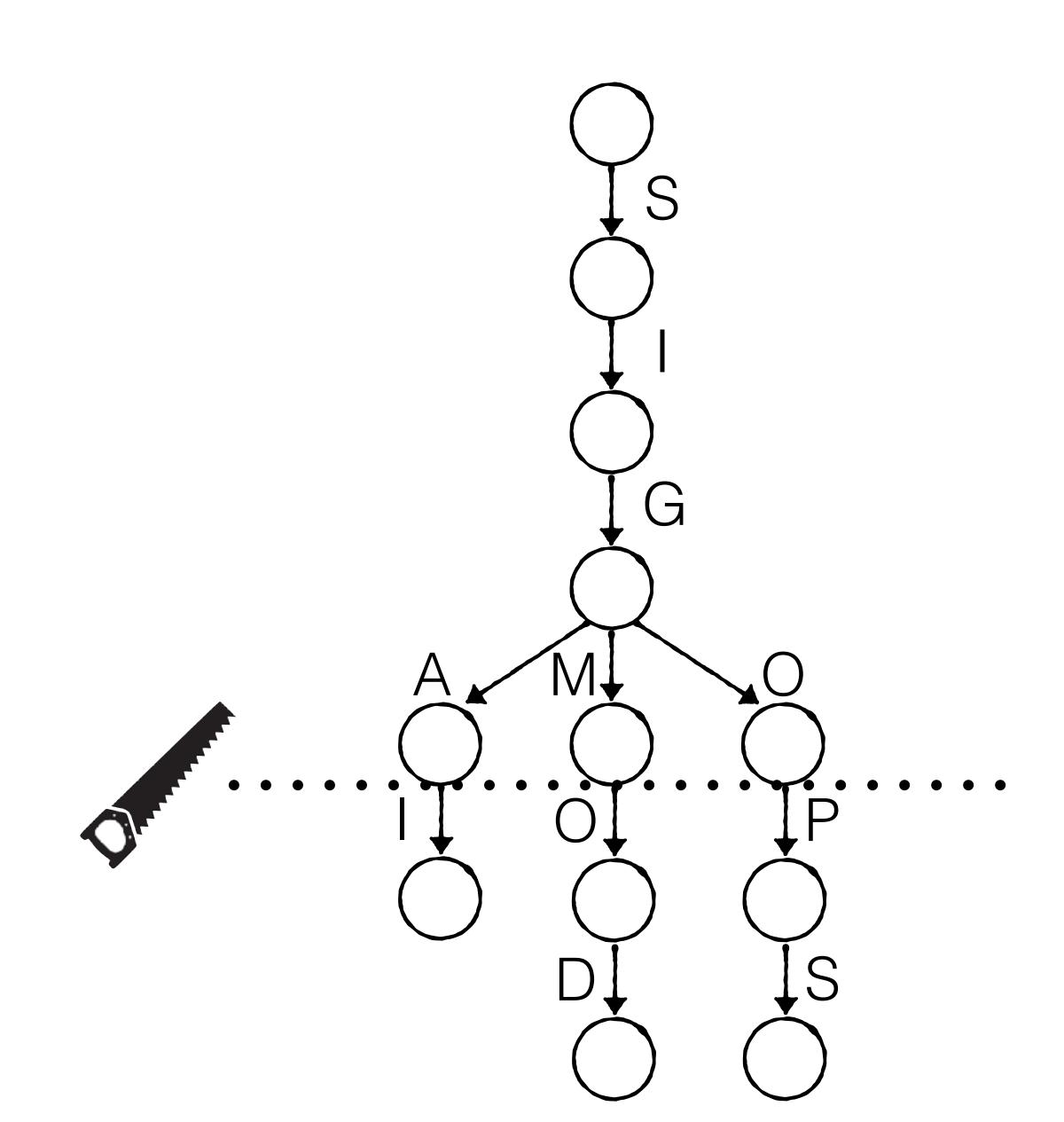


2. Pointers take space

# Stores full keys truncation

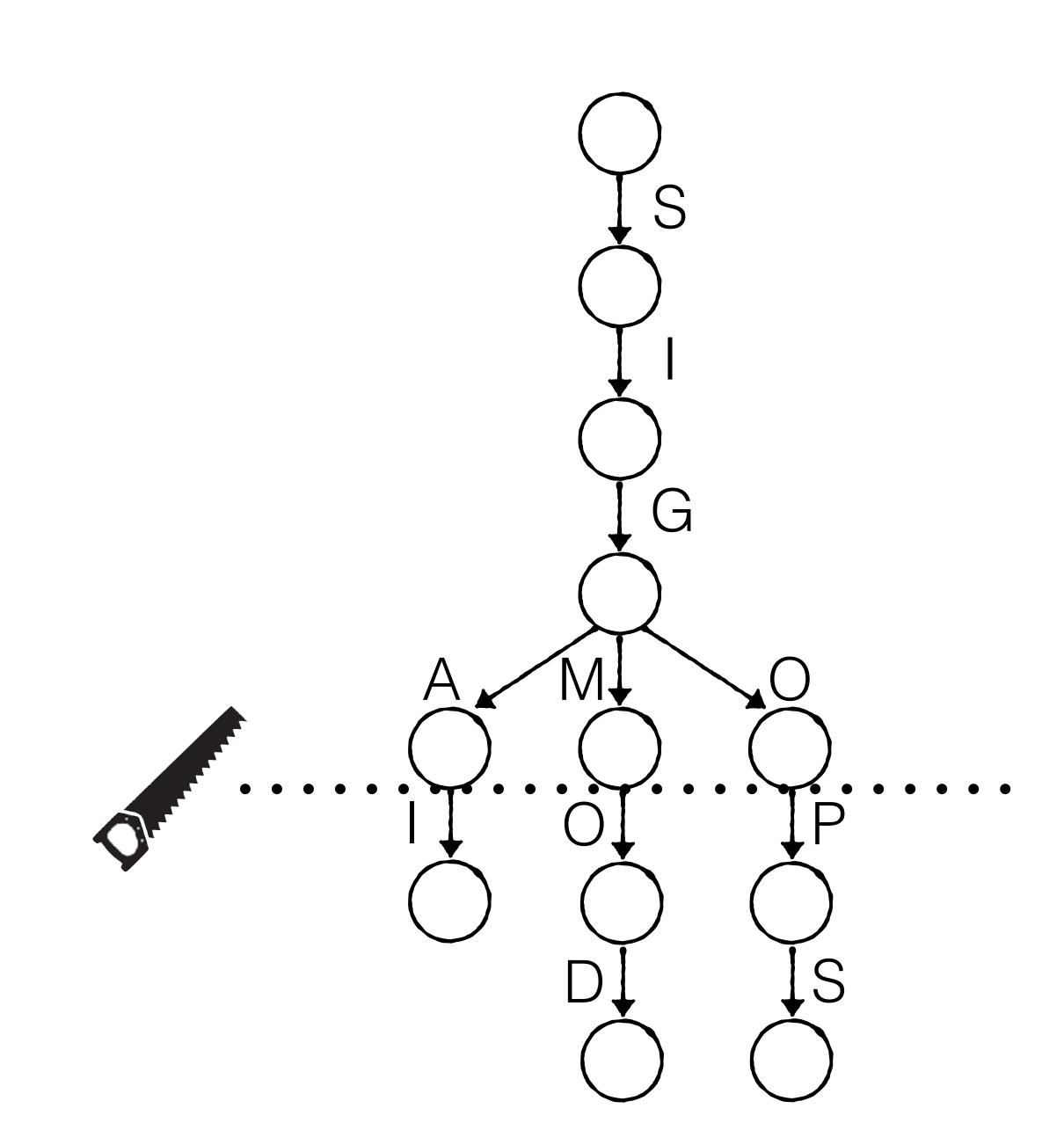


# 2. Pointers take spaceSuccinct encoding



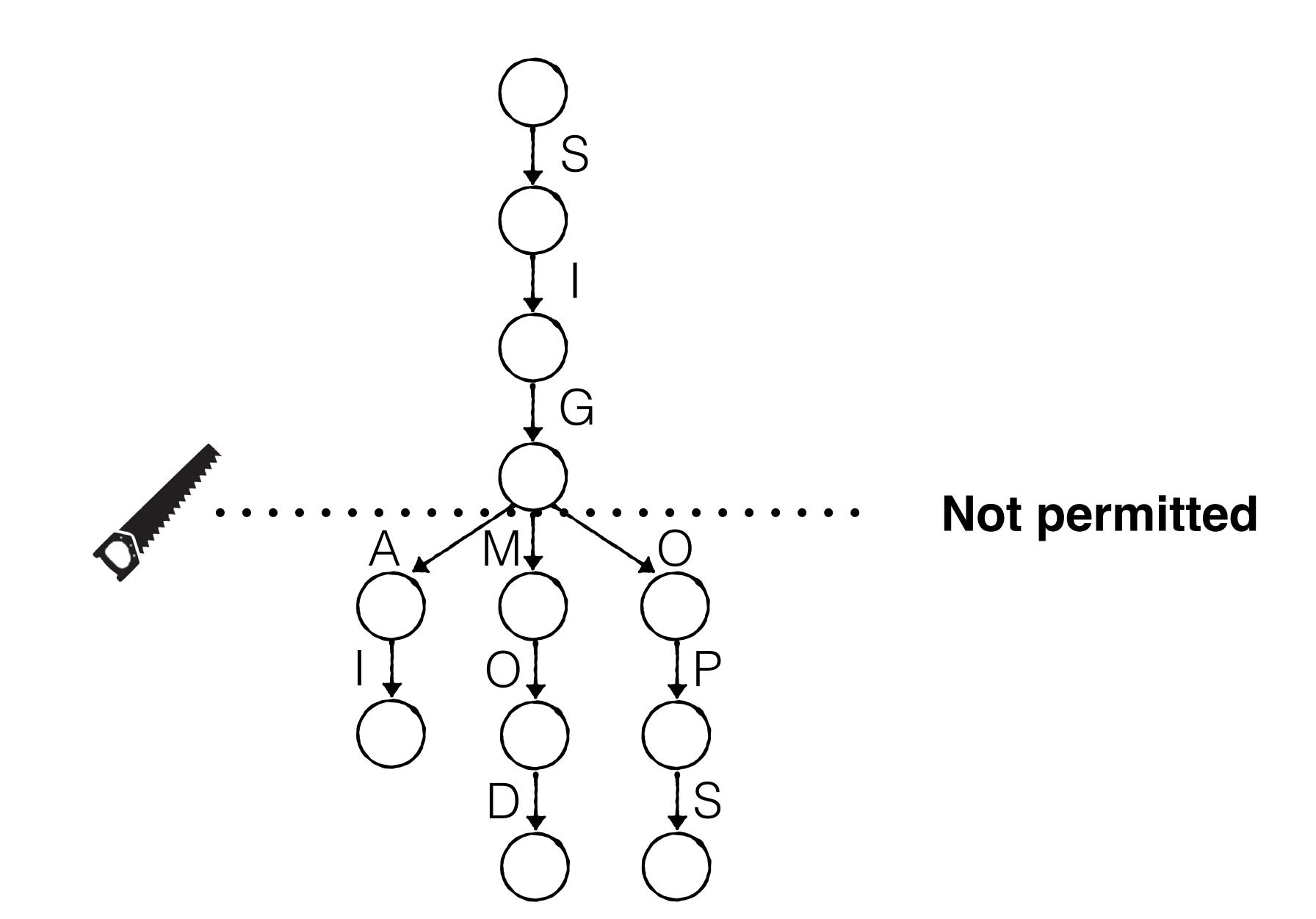
truncation

# Keep at least one unique byte for each key

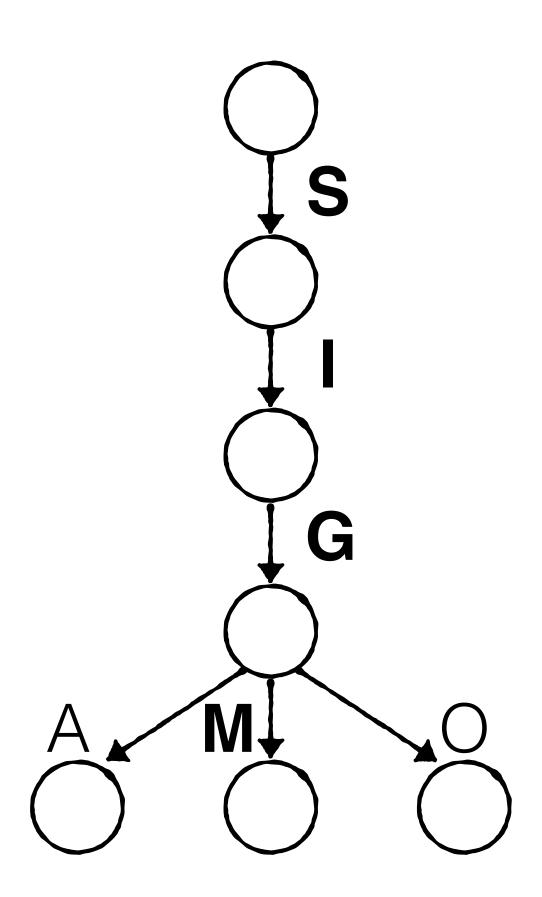


truncation

Keep at least one unique byte for each key

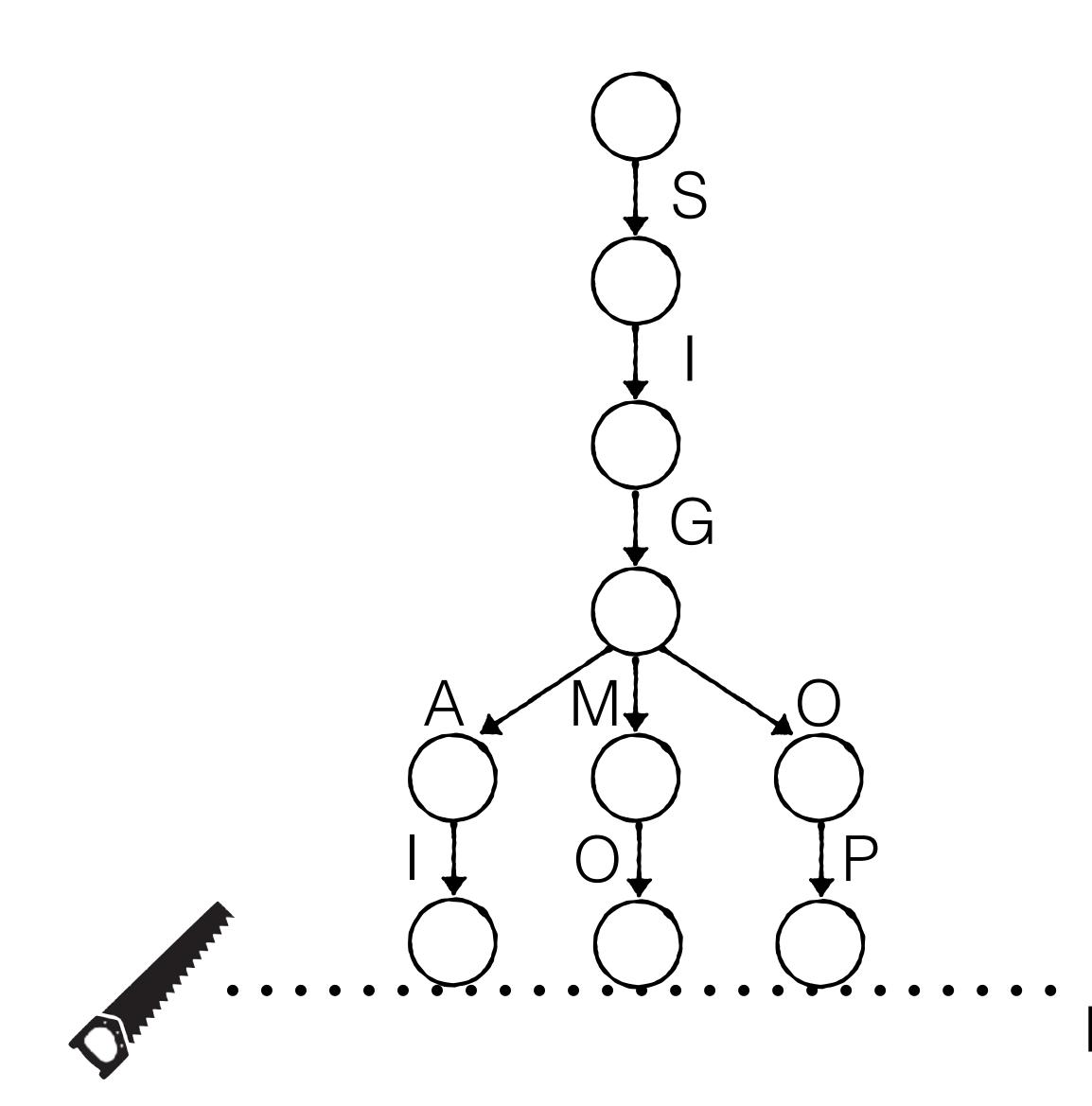


## We now get false positives (e.g., query SIGMETRICS)



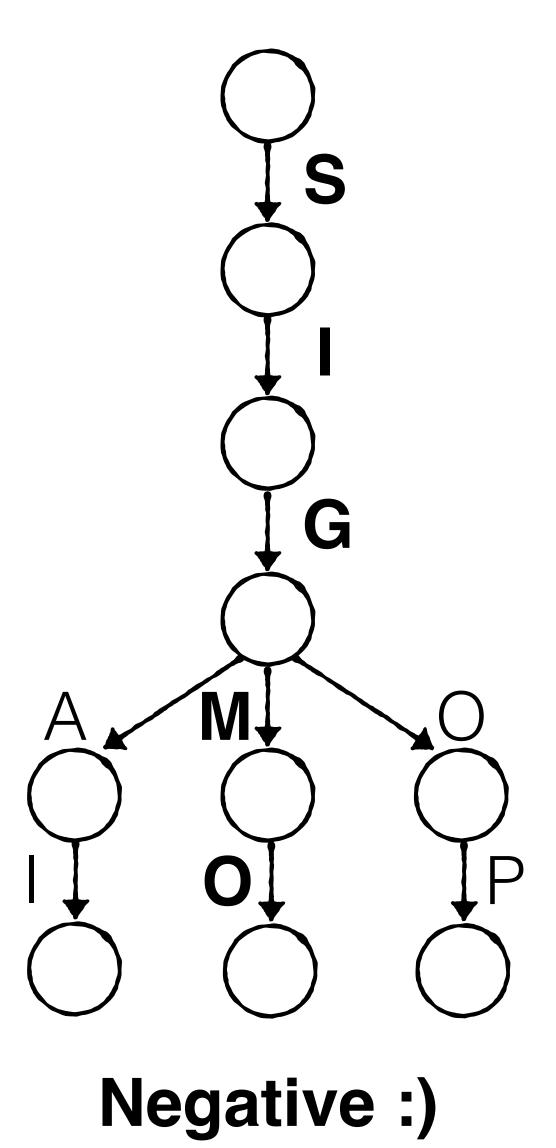
Match

We now get false positives (e.g., query SIGMETRICS)

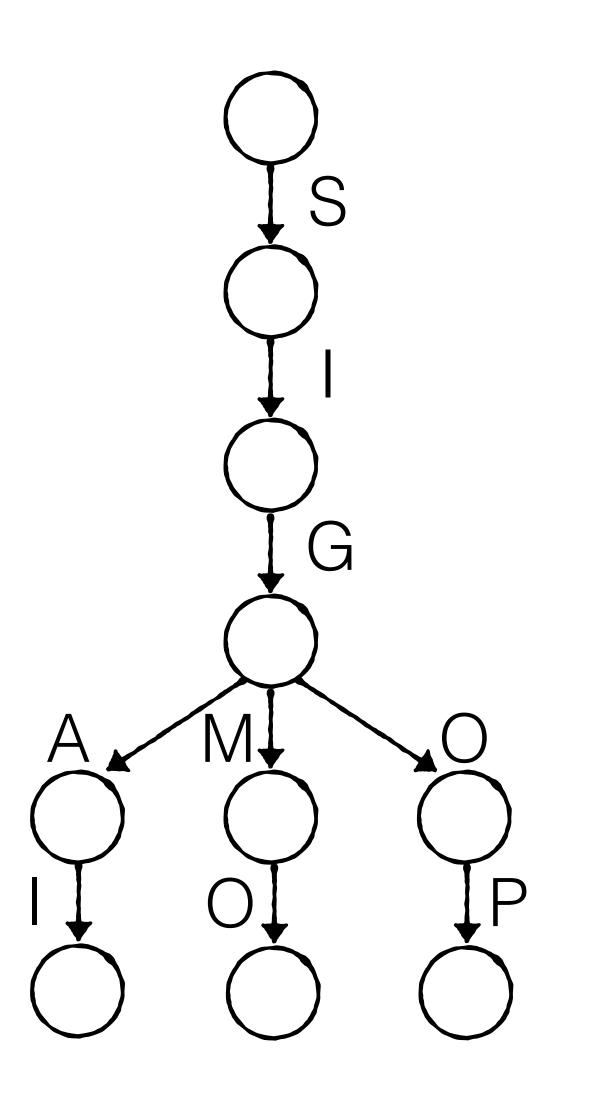


Can truncate later - better FPR, worse space

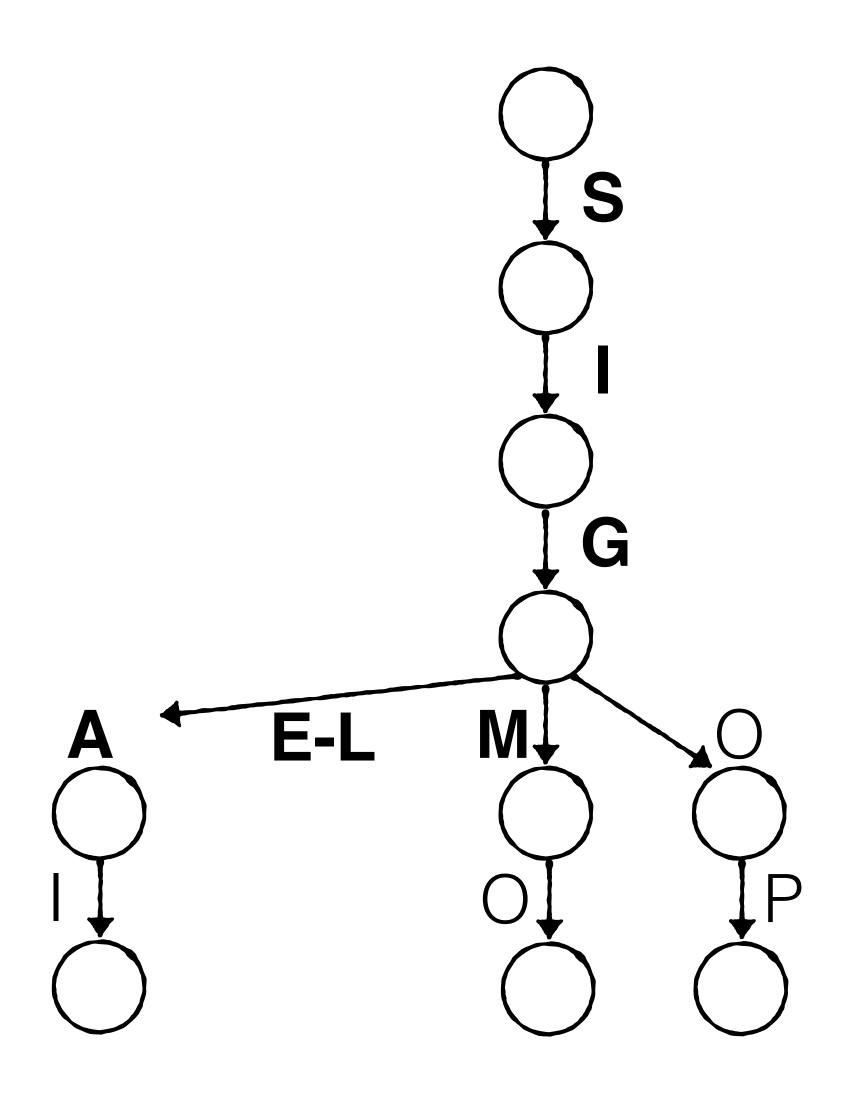
We now get false positives (e.g., query SIGMETRICS)



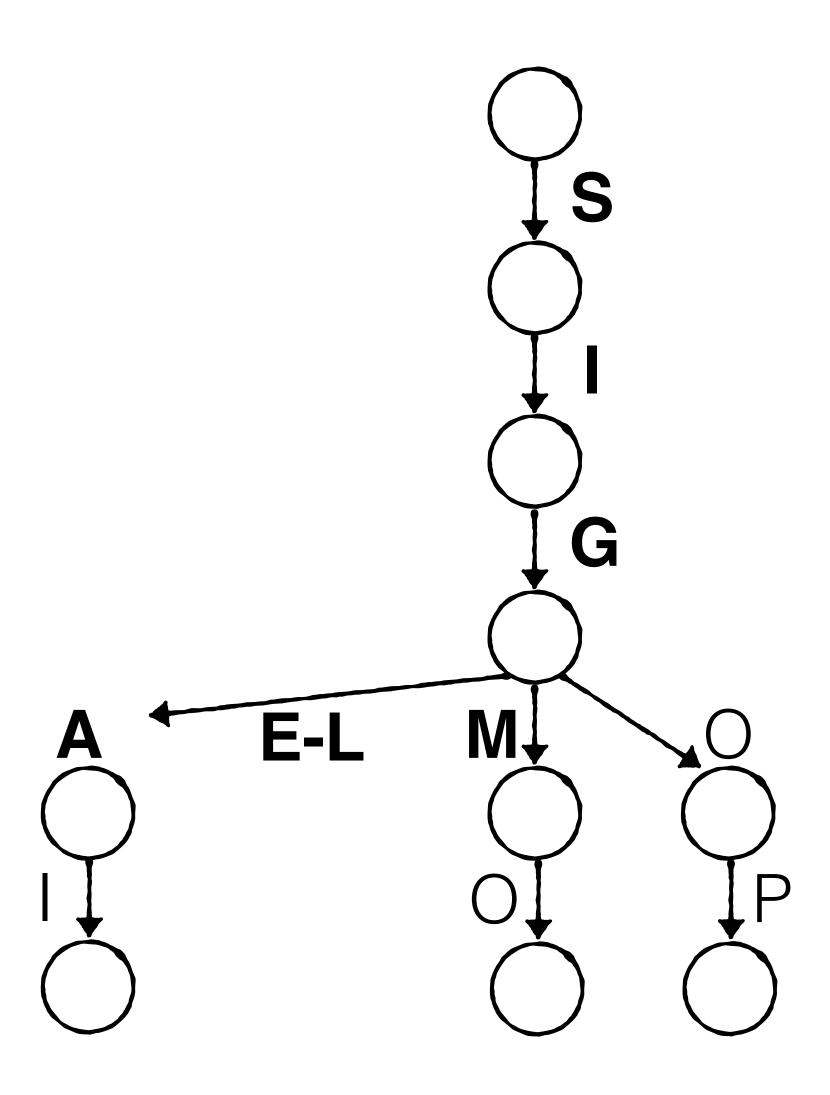
## Range query: SIGE - SIGL



Range query: SIGE - SIGL

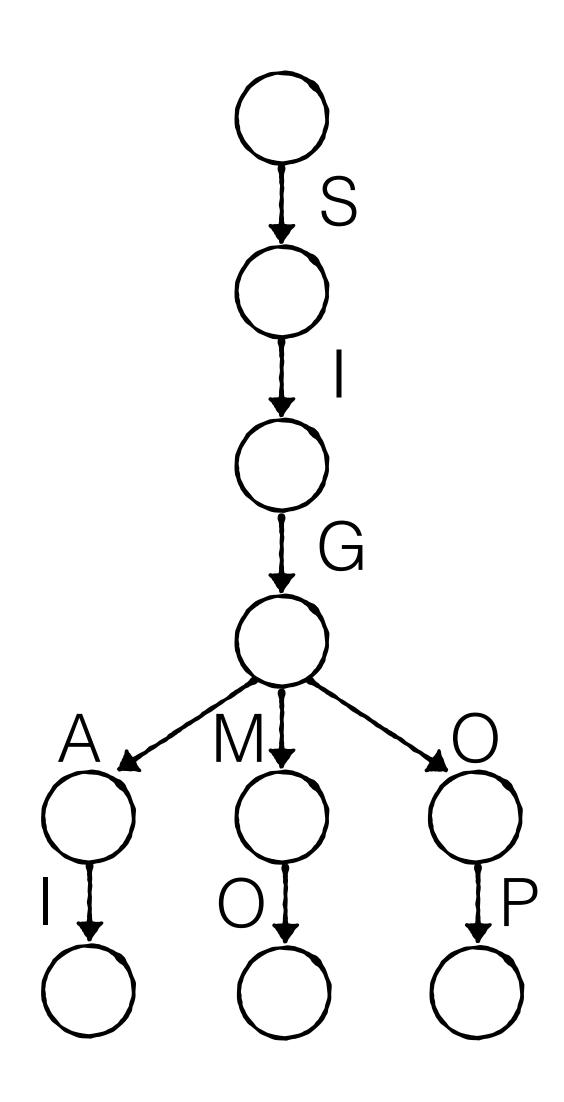


# Range query: SIGE - SIGL

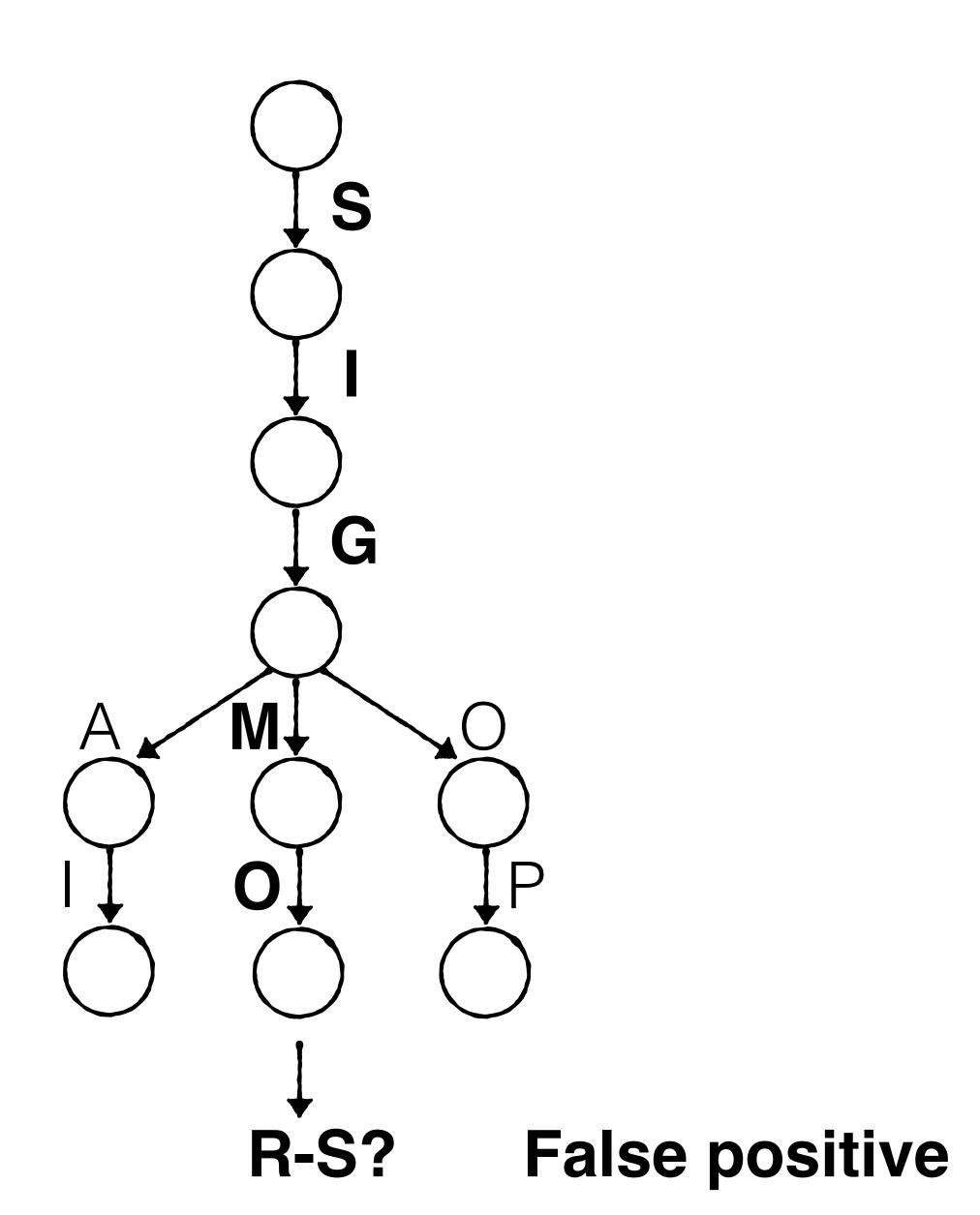


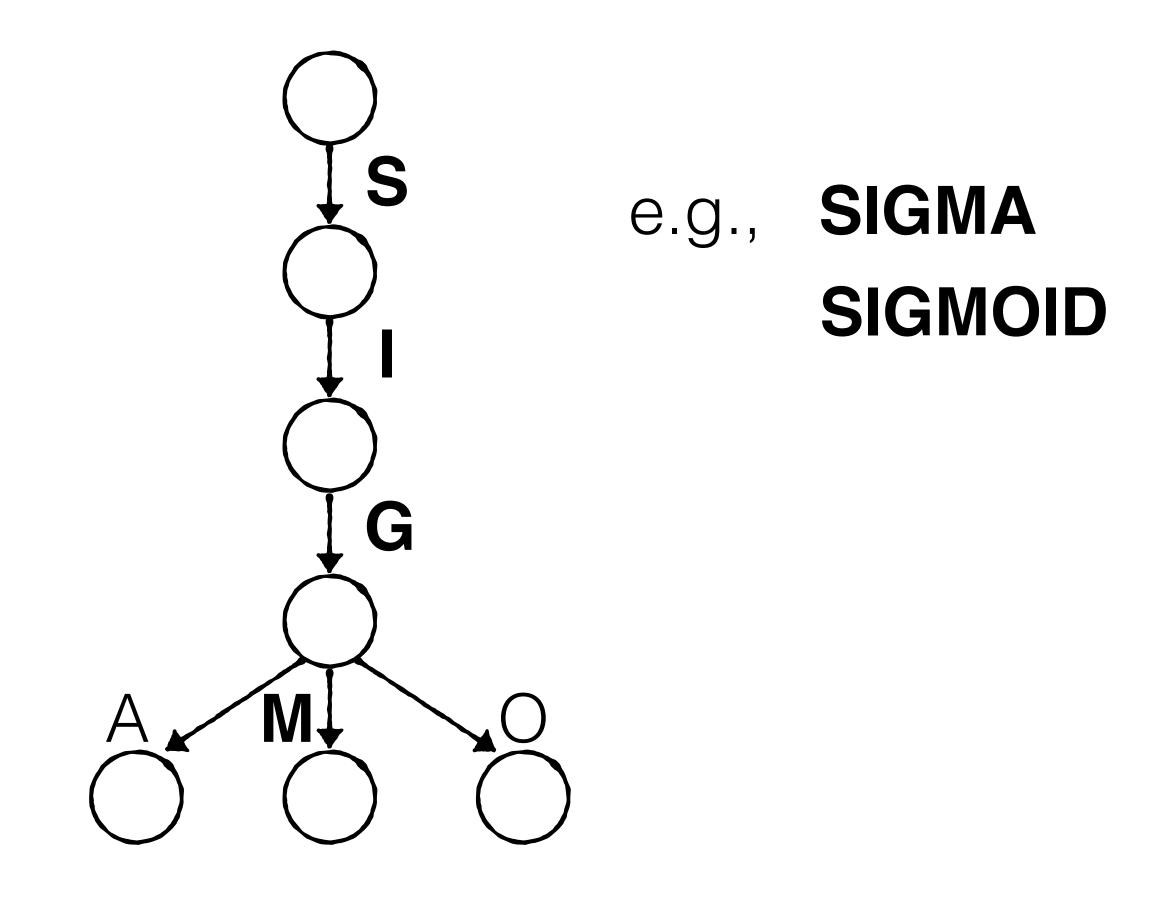
**Negative:)** 

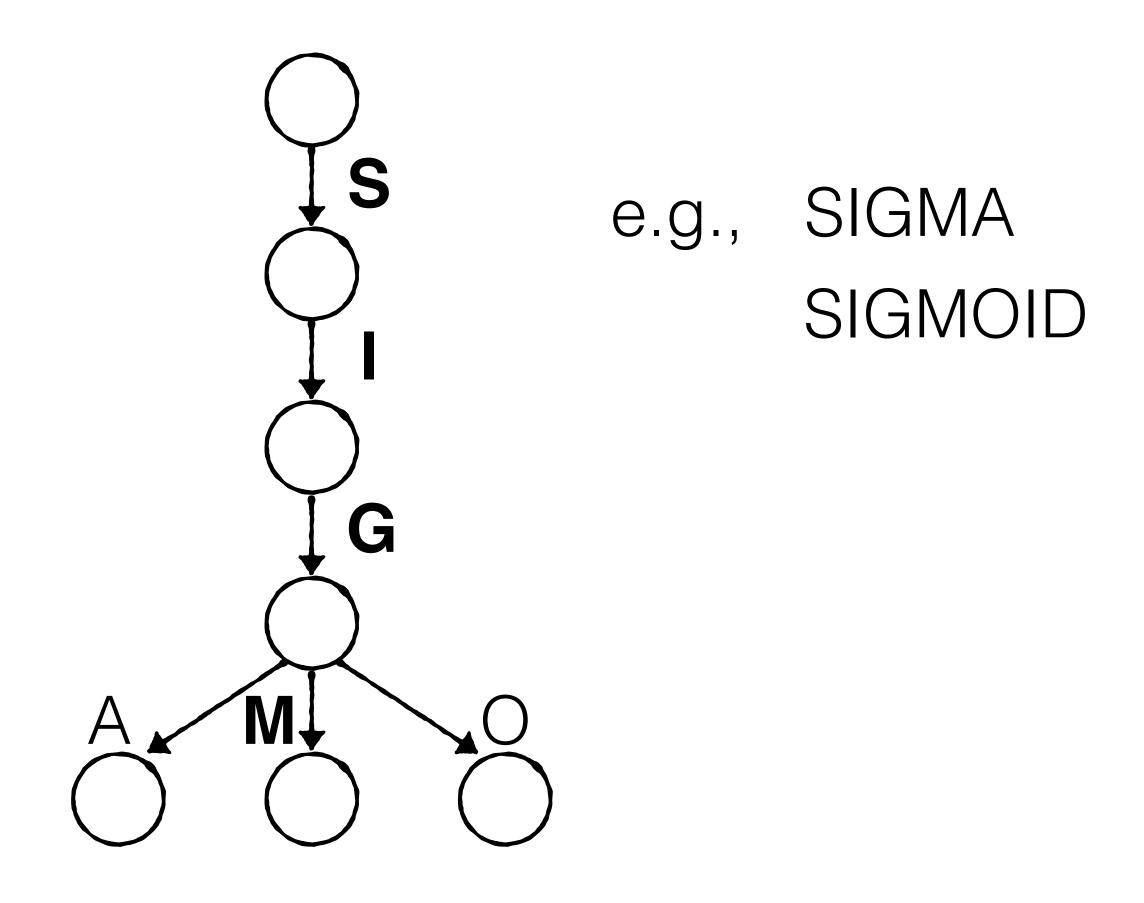
# Range query: SIGMOR - SIGMOS



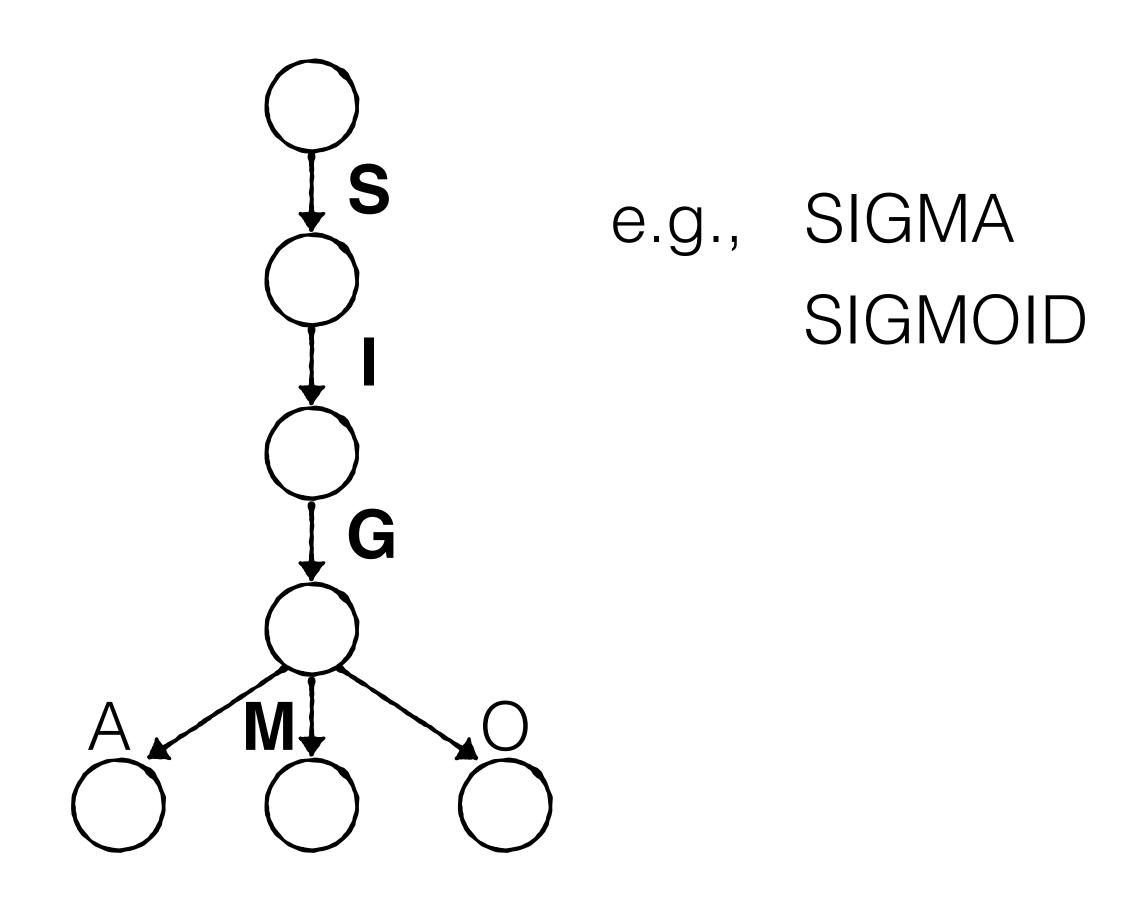
## Range query: SIGMOR - SIGMOS





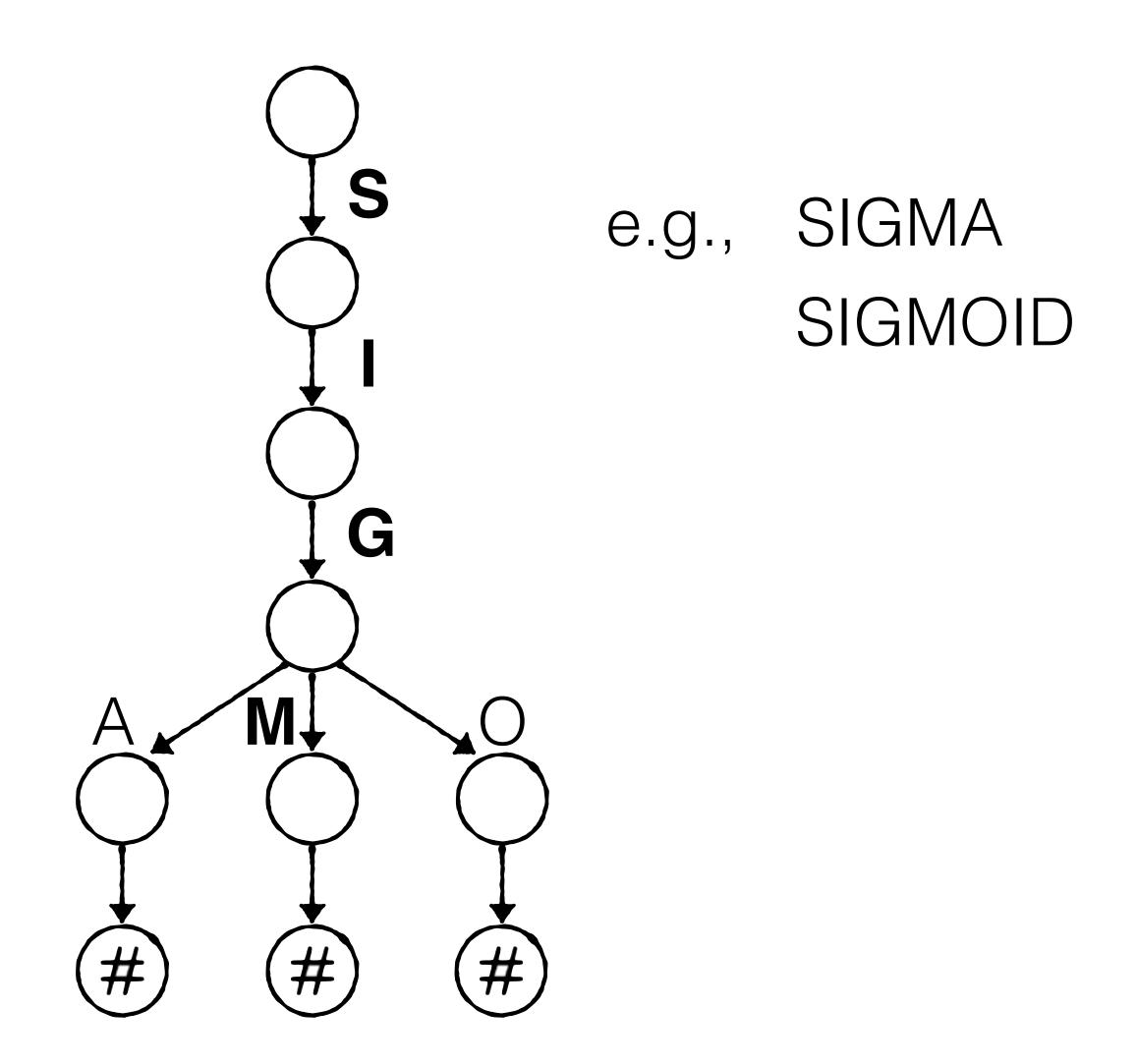


**Correlated Workload** 



Correlated Workload

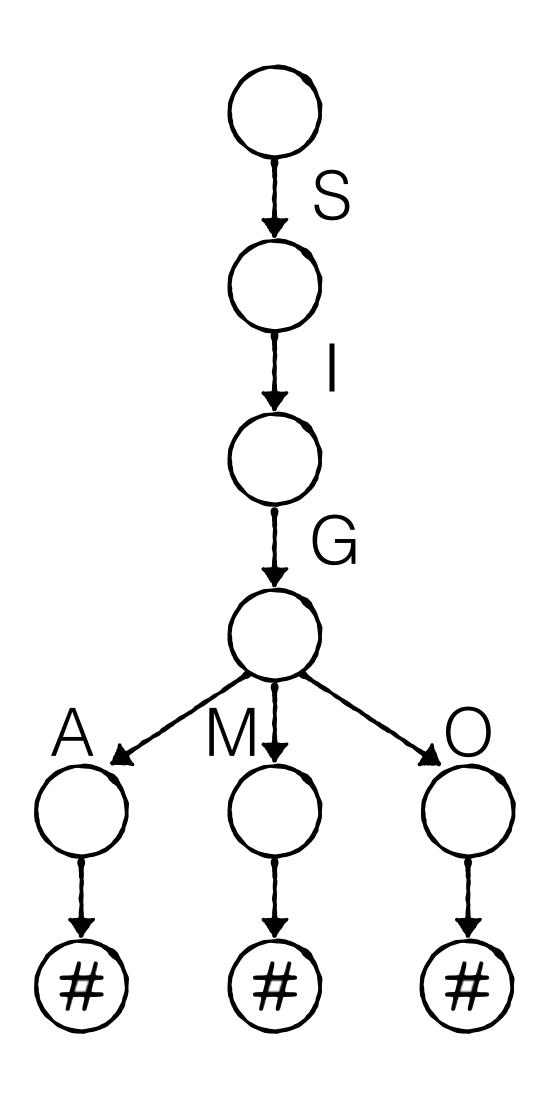
Can we alleviate this problem for point queries?



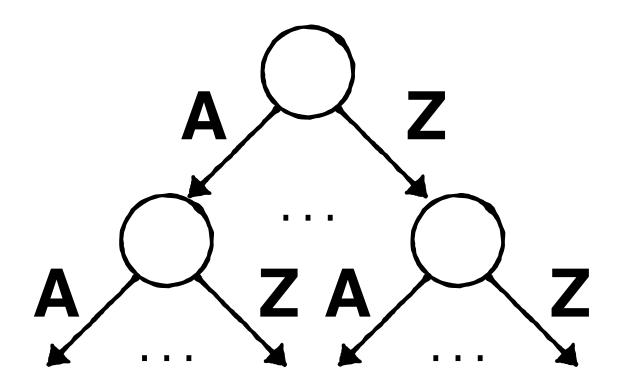
Can we alleviate this problem for point queries?

Add 1 byte hash

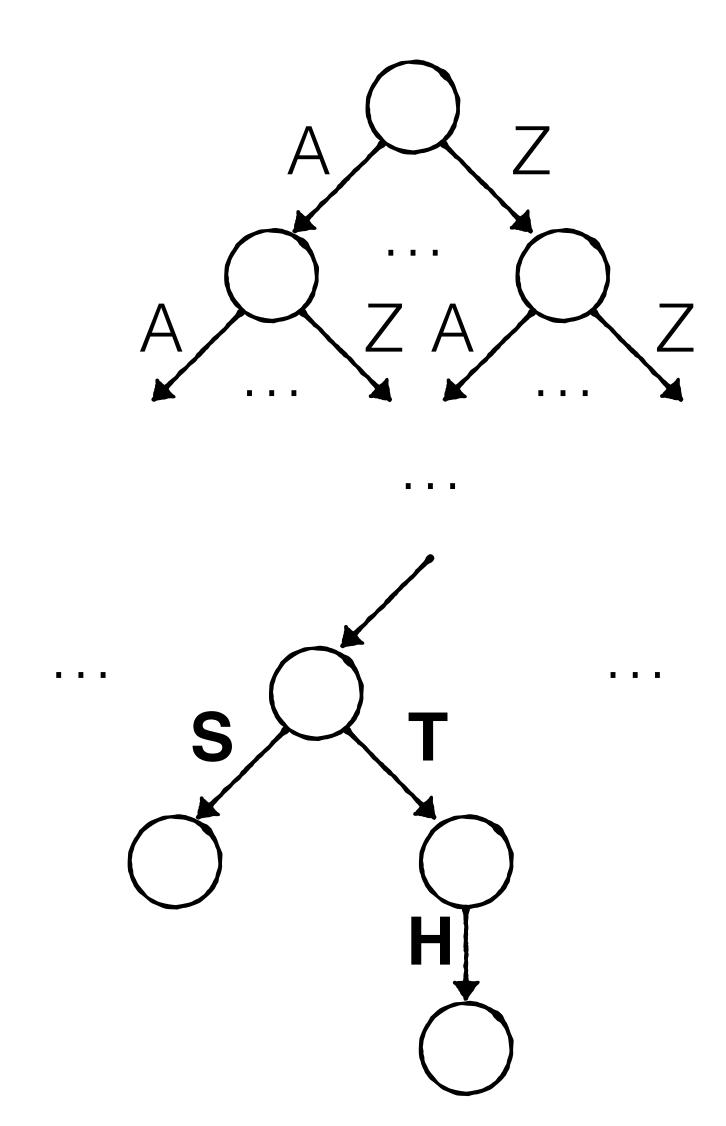
## How to encode the trie succinctly?



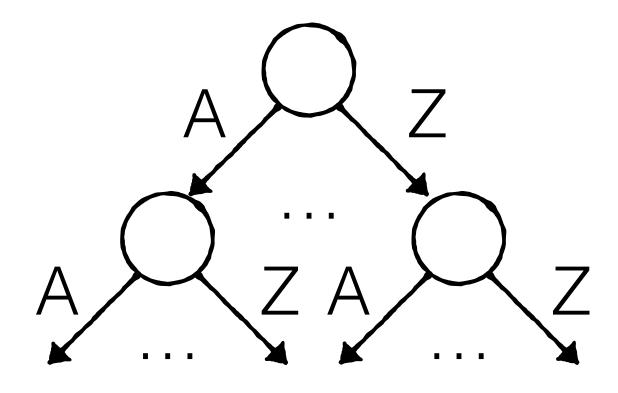
## Typically, the upper levels of the trie are more full

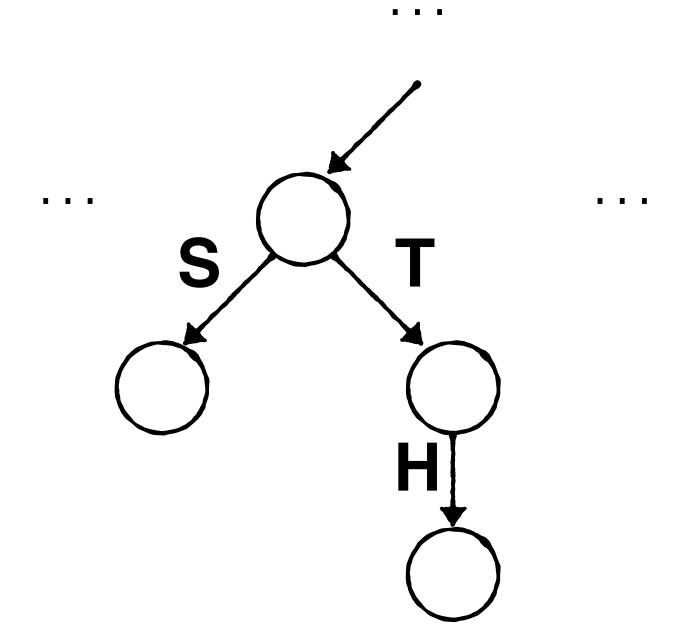


## Typically, the upper levels of the trie are more full

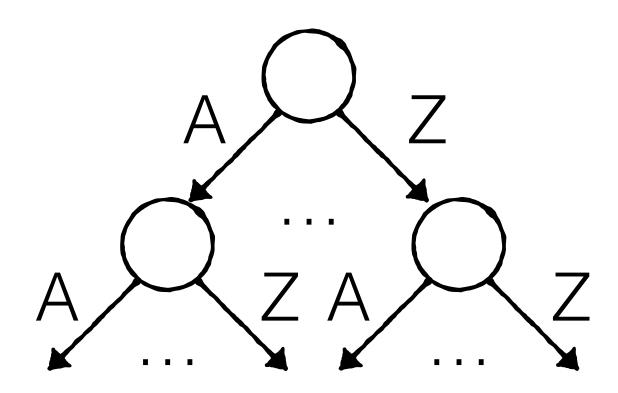


# More queries go through base layers

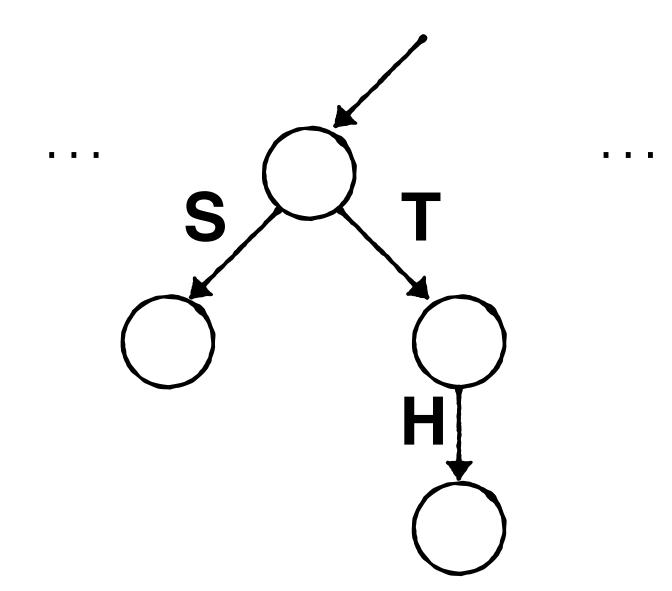




More queries go through base layers

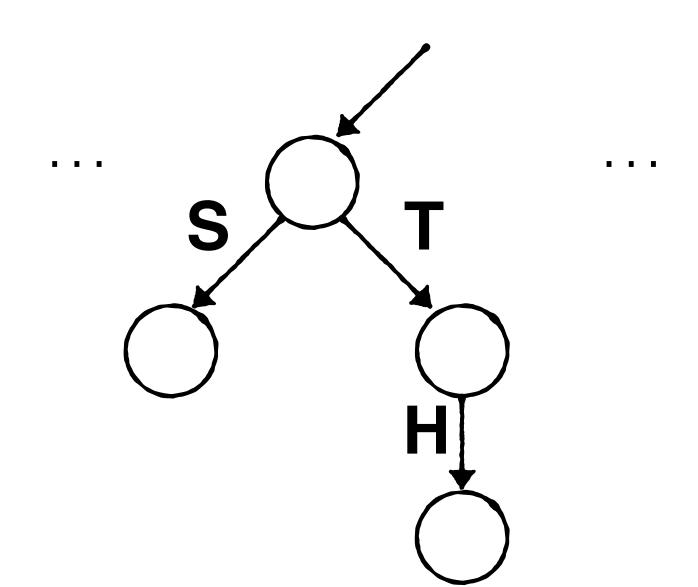


A leaf is reached on avg. by fewer queries

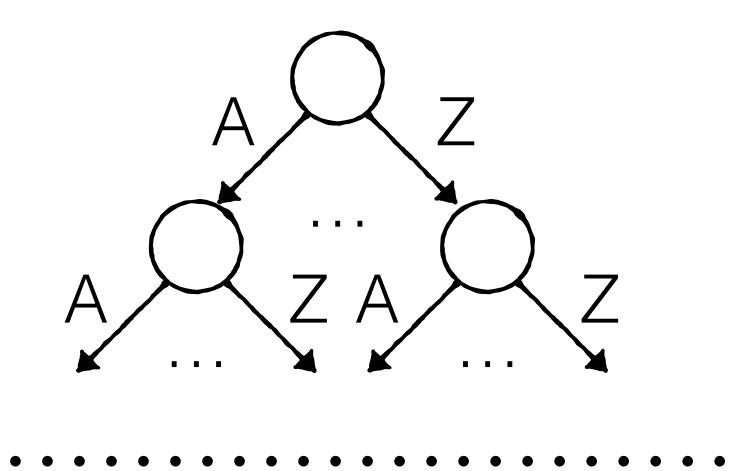


More queries go through base layers A Z Z Z

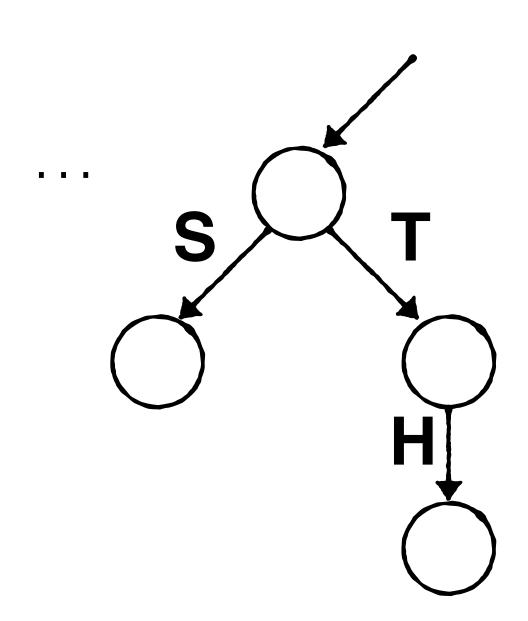
A leaf is reached on avg. by fewer queries



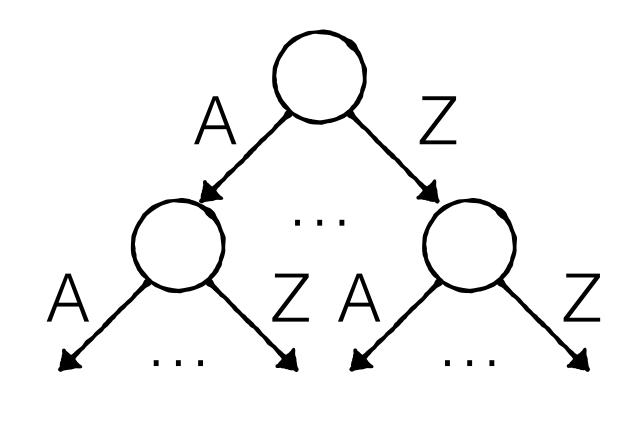
# Exponentially more nodes as we move down



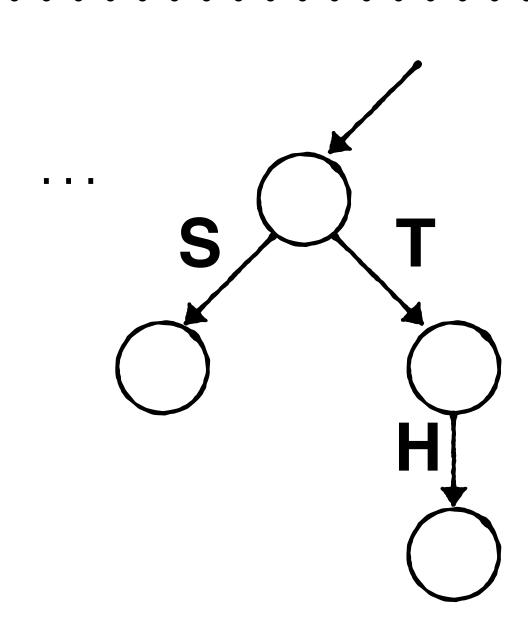
# **Optimize for speed**



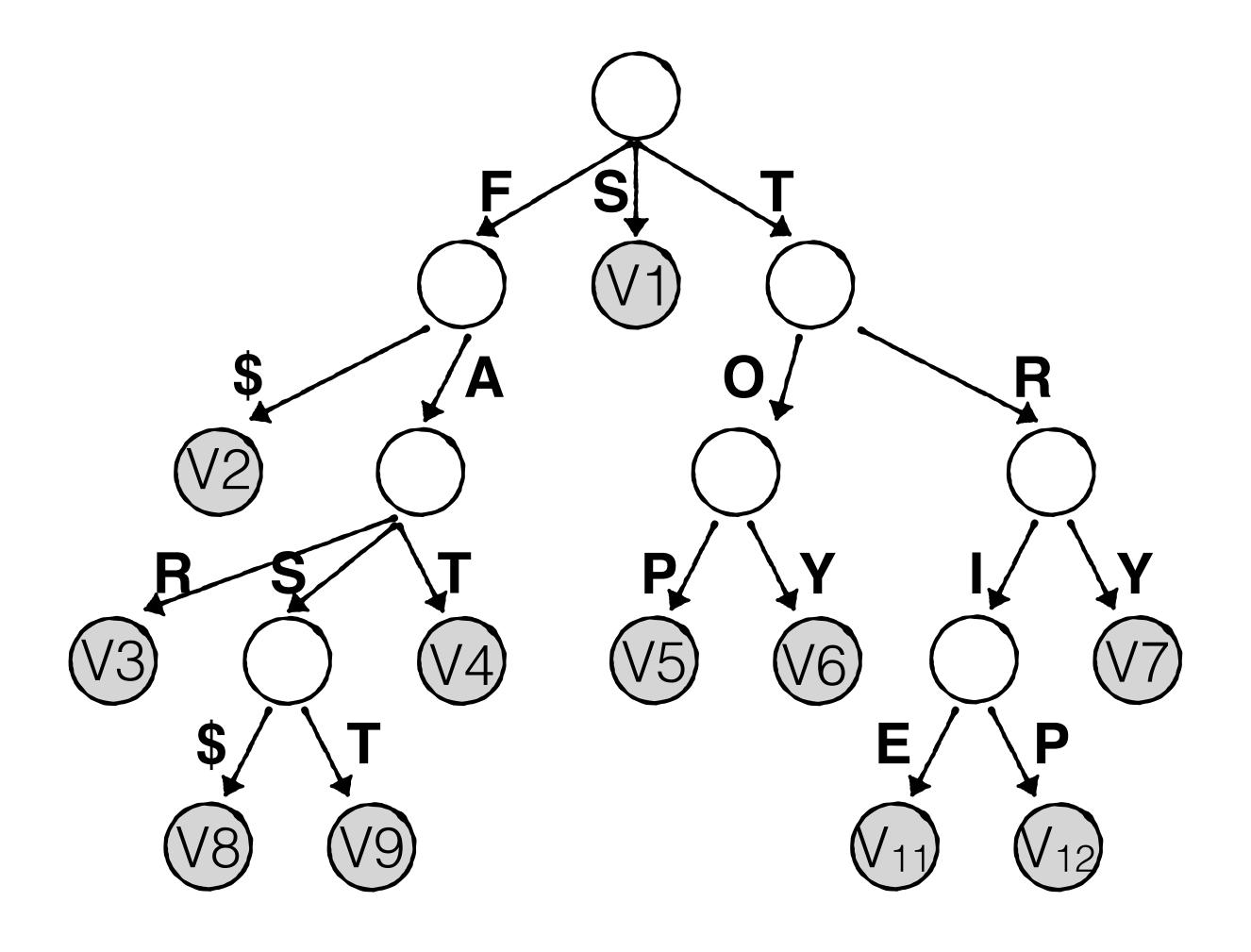
**Optimize for space** 



Optimize for speed 1/64

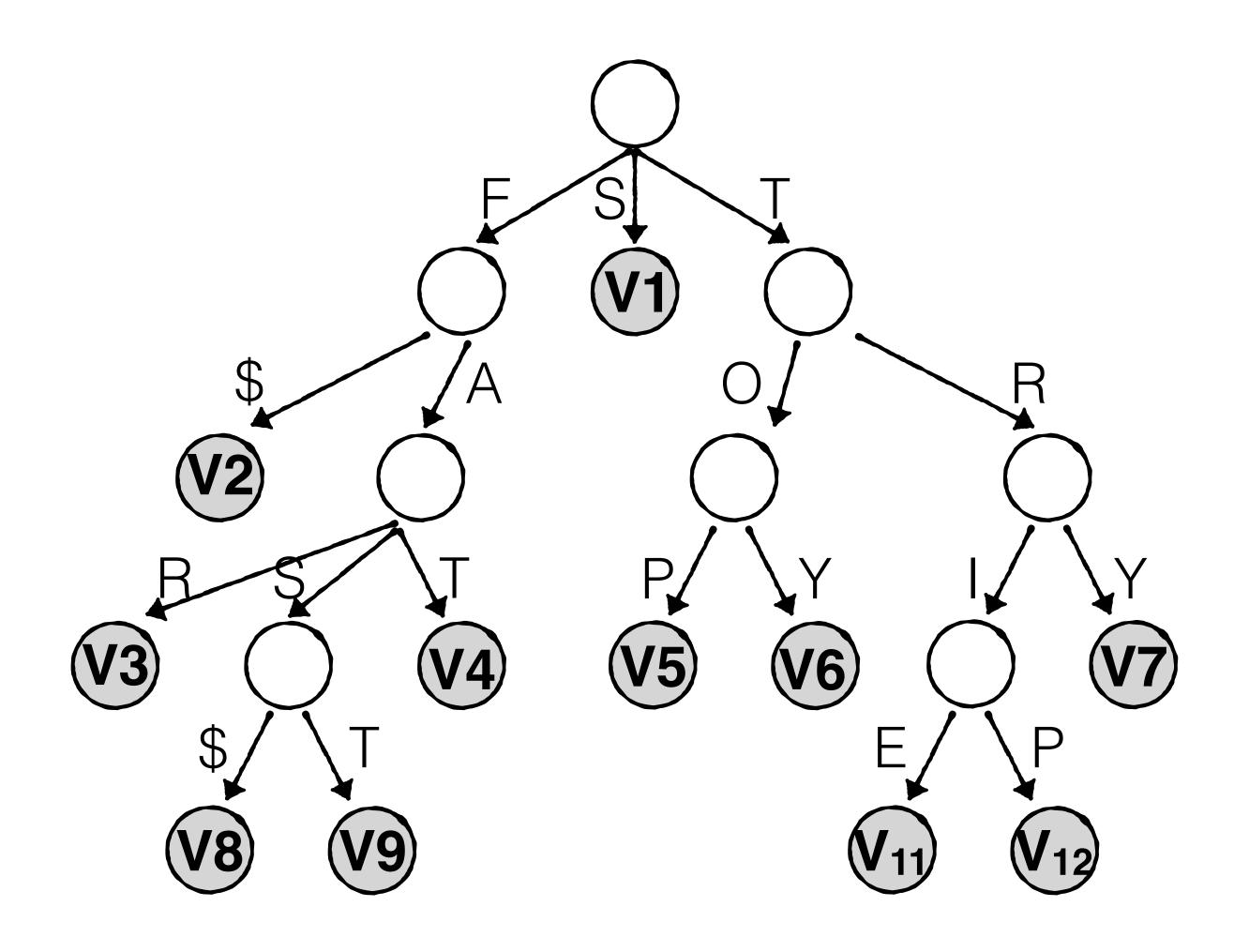


Optimize for space 63/64



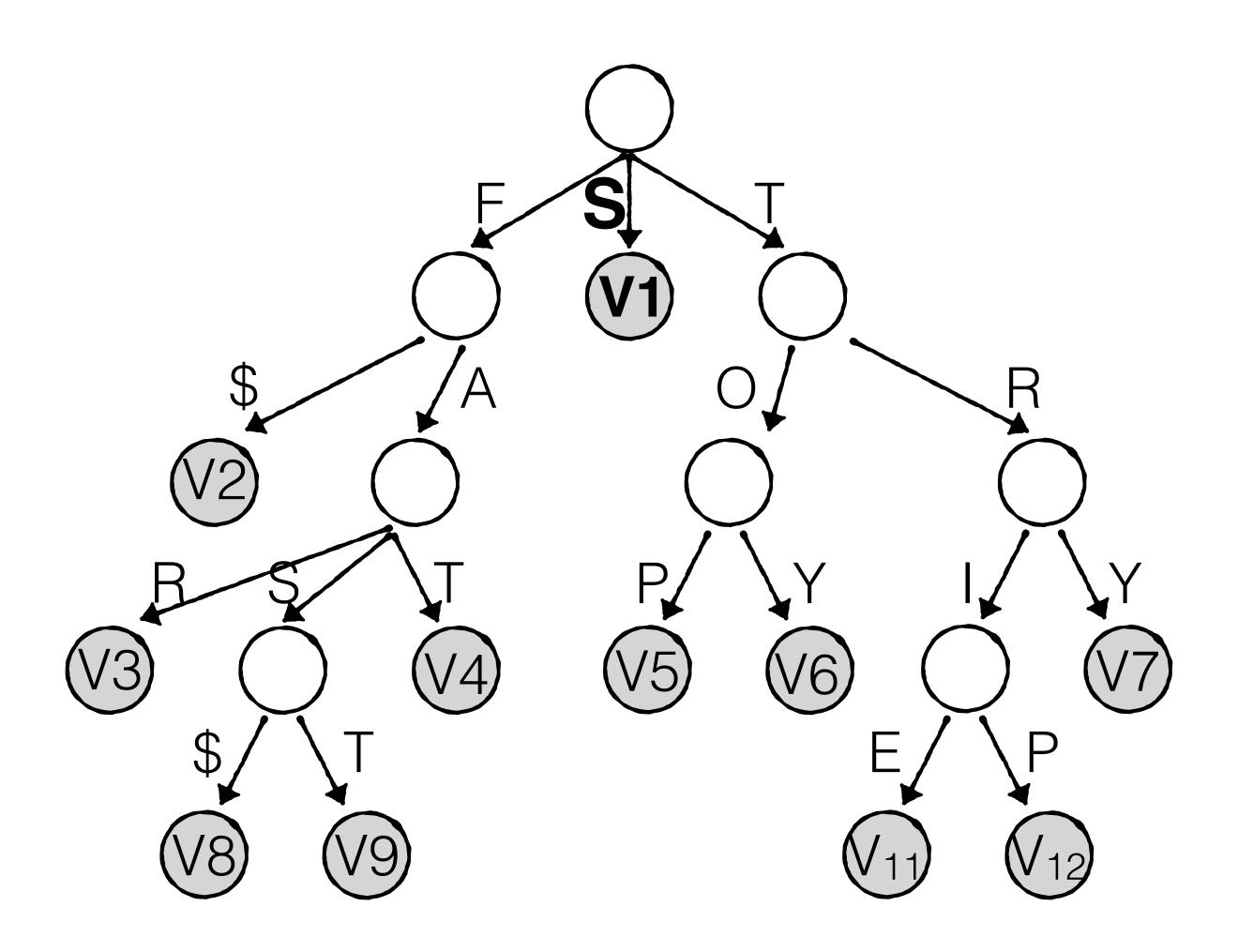
F, FAR, FAS, FAST, FAT, S, TOP, TOY, TRIE, TRIP, TRY

## Each key leads to leaf payload (e.g., pointer and/or hash)

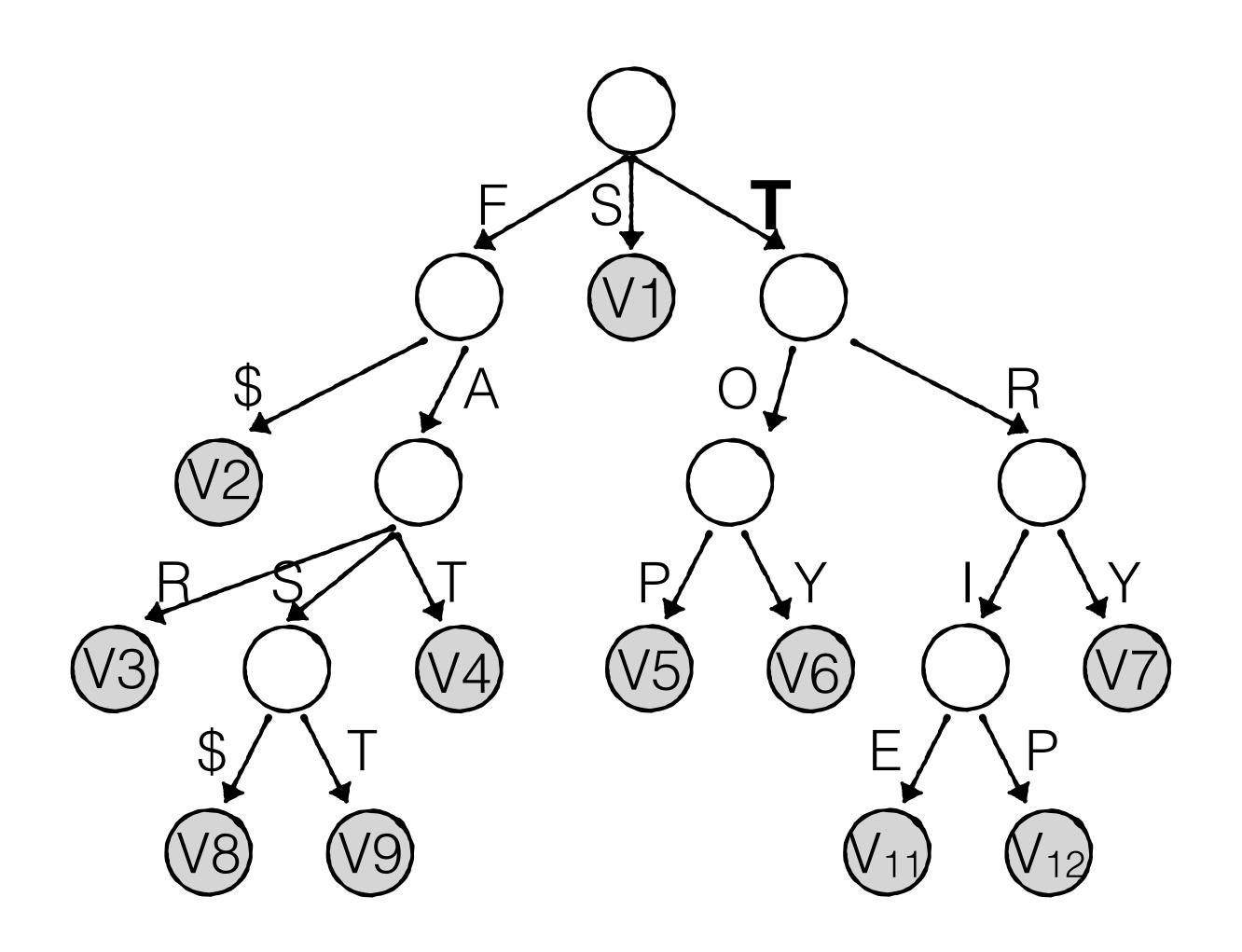


F, FAR, FAS, FAST, FAT, S, TOP, TOY, TRIE, TRIP, TRY

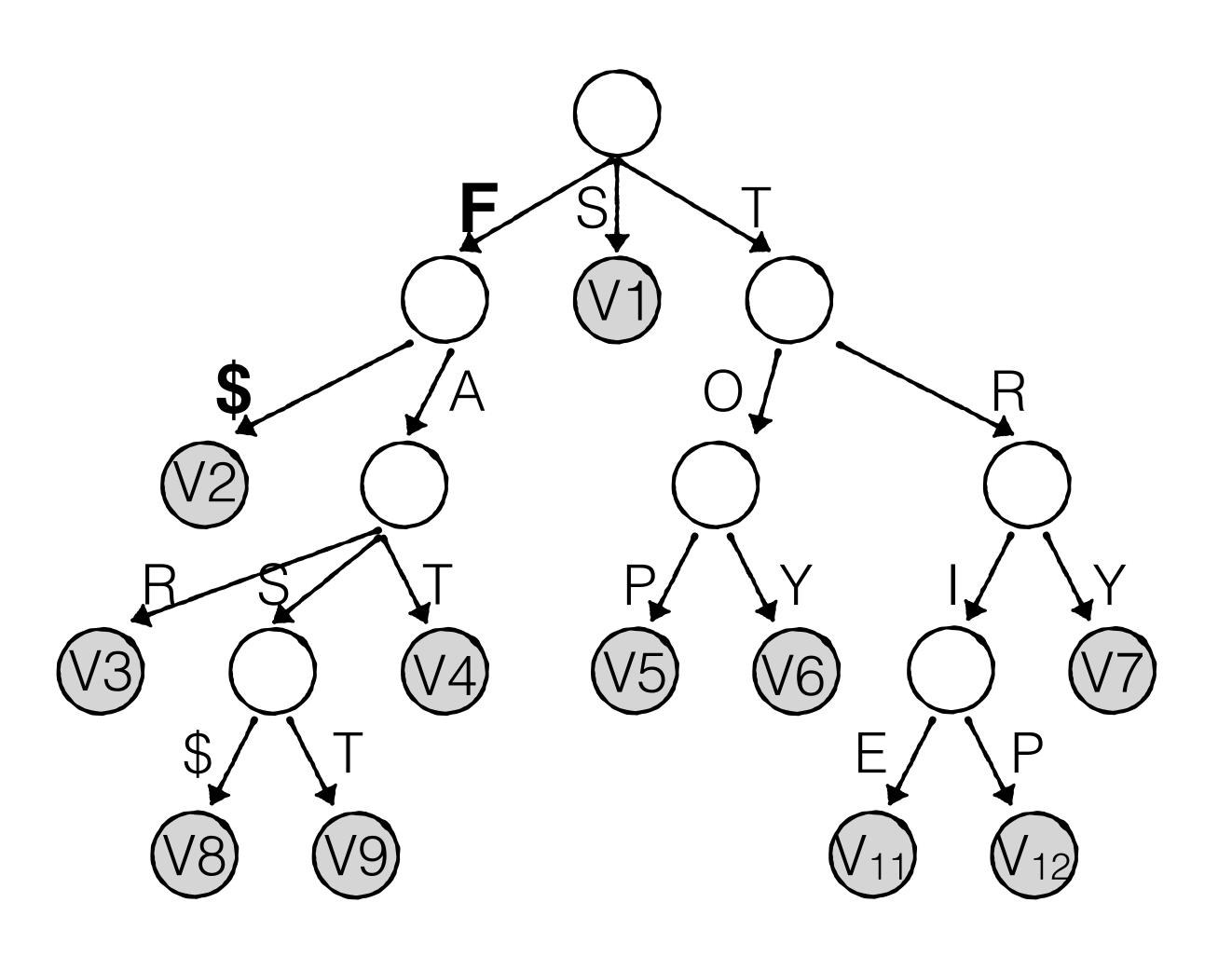
S leads to leaf: full key not prefix of any other key



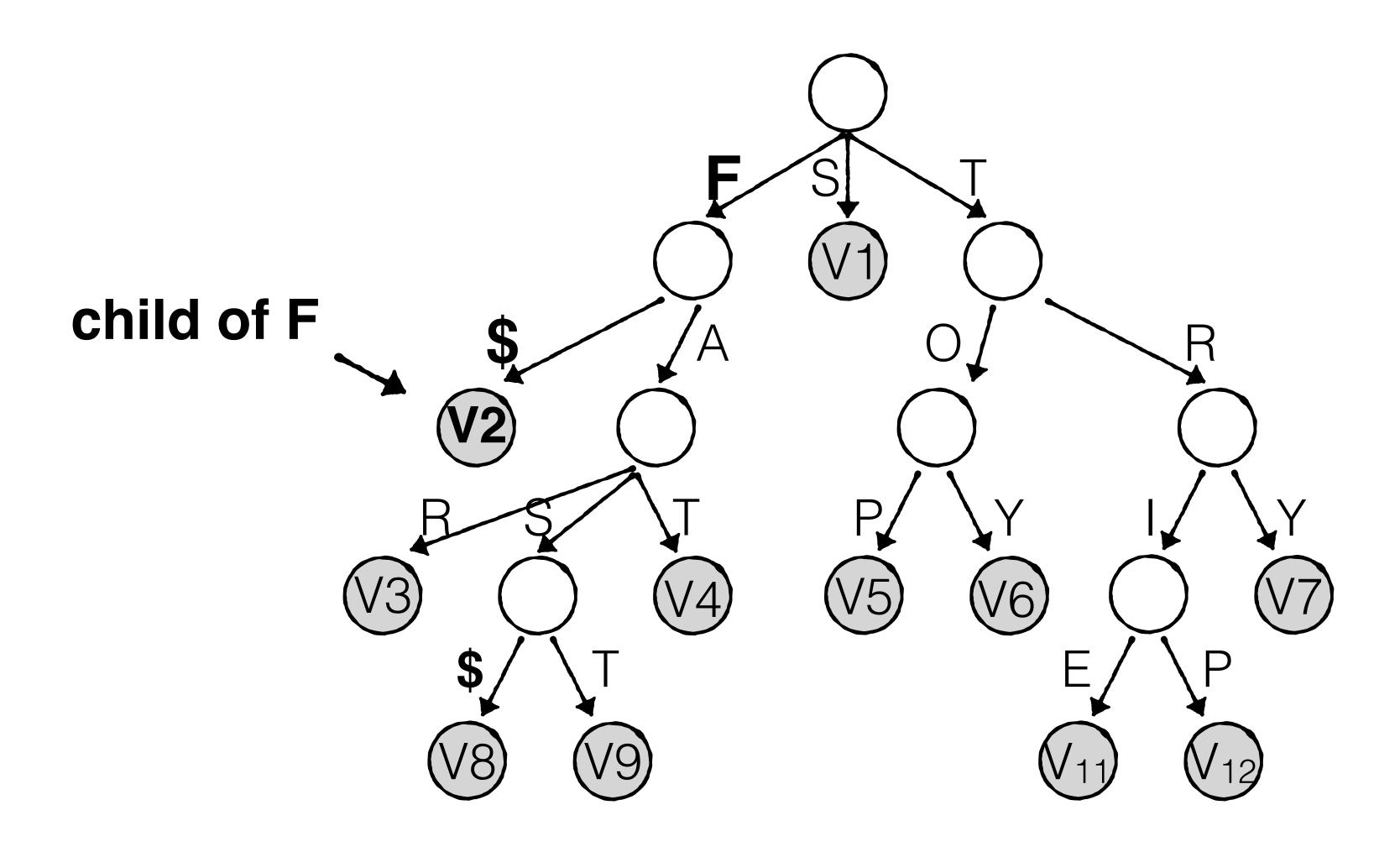
T leads to internal node: not full key prefix of one or more other keys



F leads to internal node with \$ child: full key prefix of one or more other keys

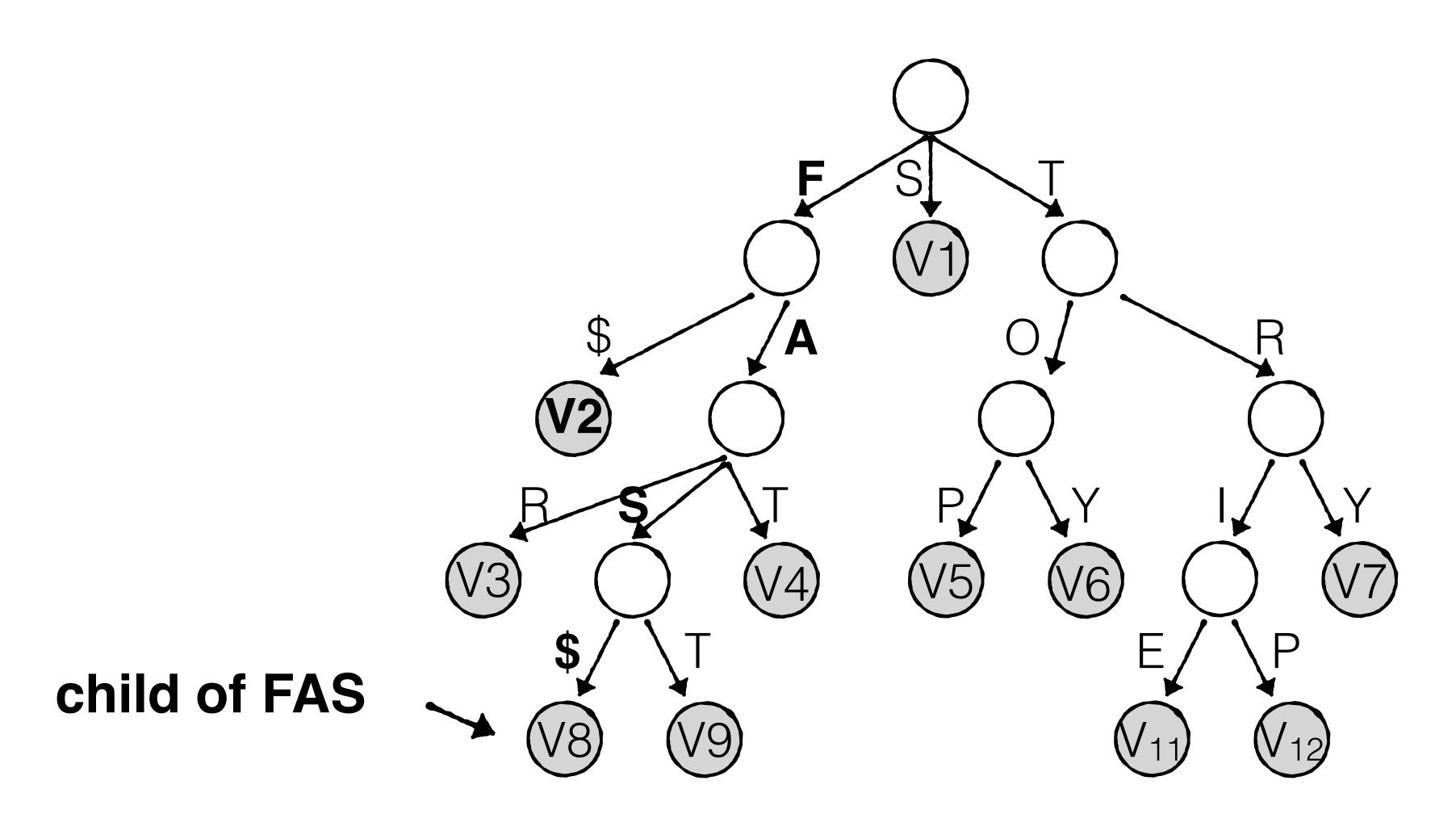


F leads to internal node with \$ child: full key prefix of one or more other keys

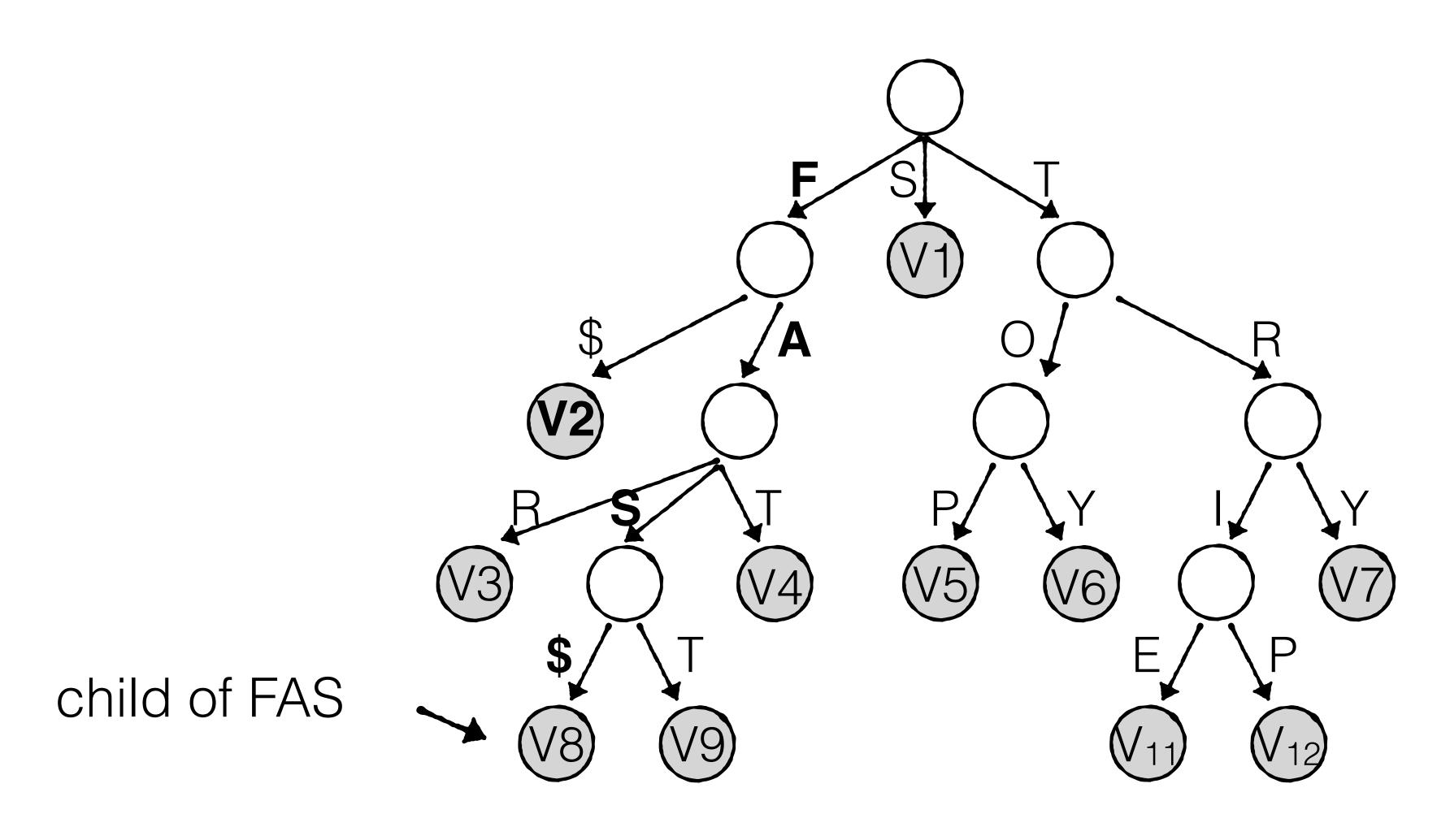


F, FAR, FAS, FAST, FAT, S, TOP, TOY, TRIE, TRIP, TRY

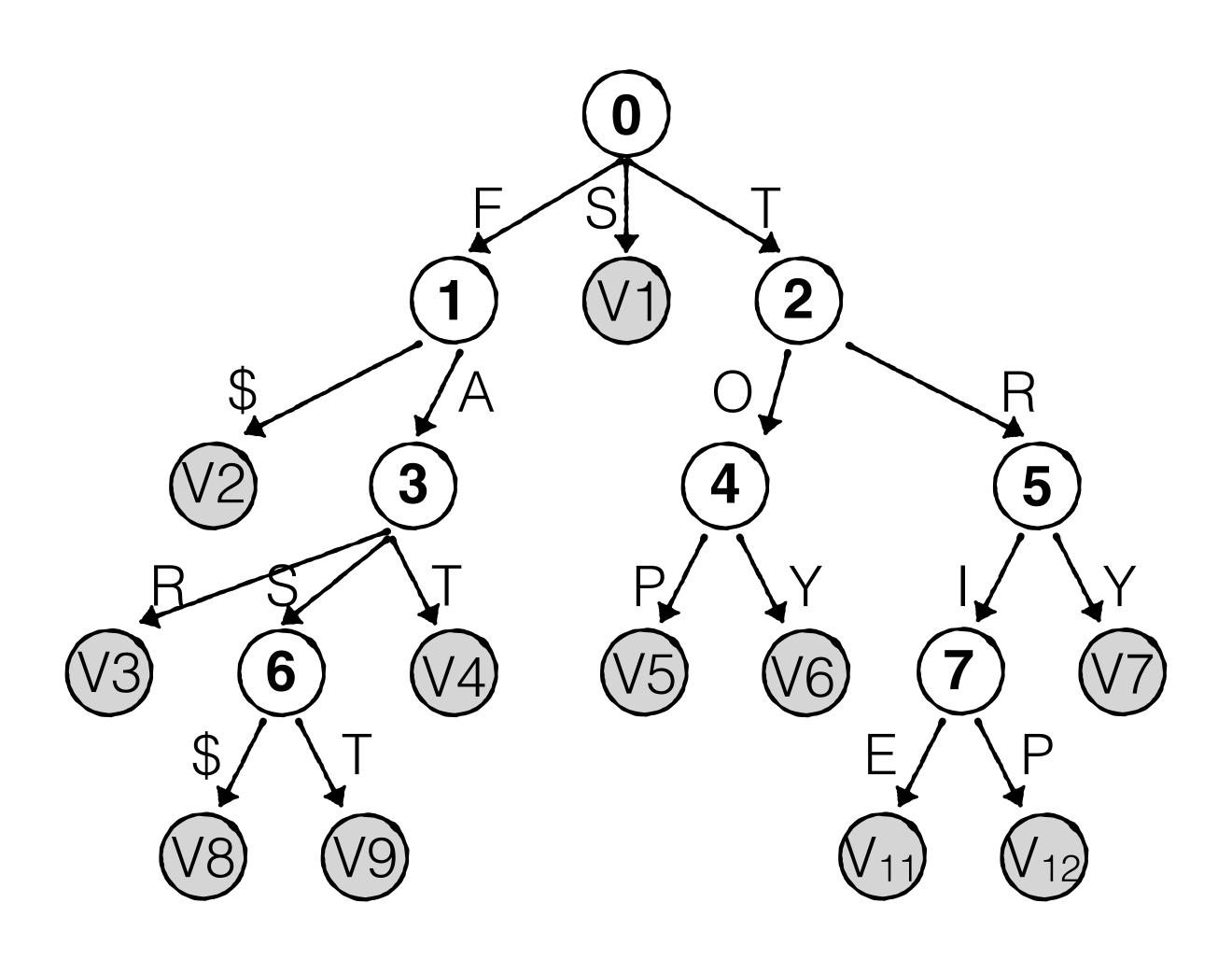
**FAS** leads to internal node with \$ child: full key prefix of one or more other keys

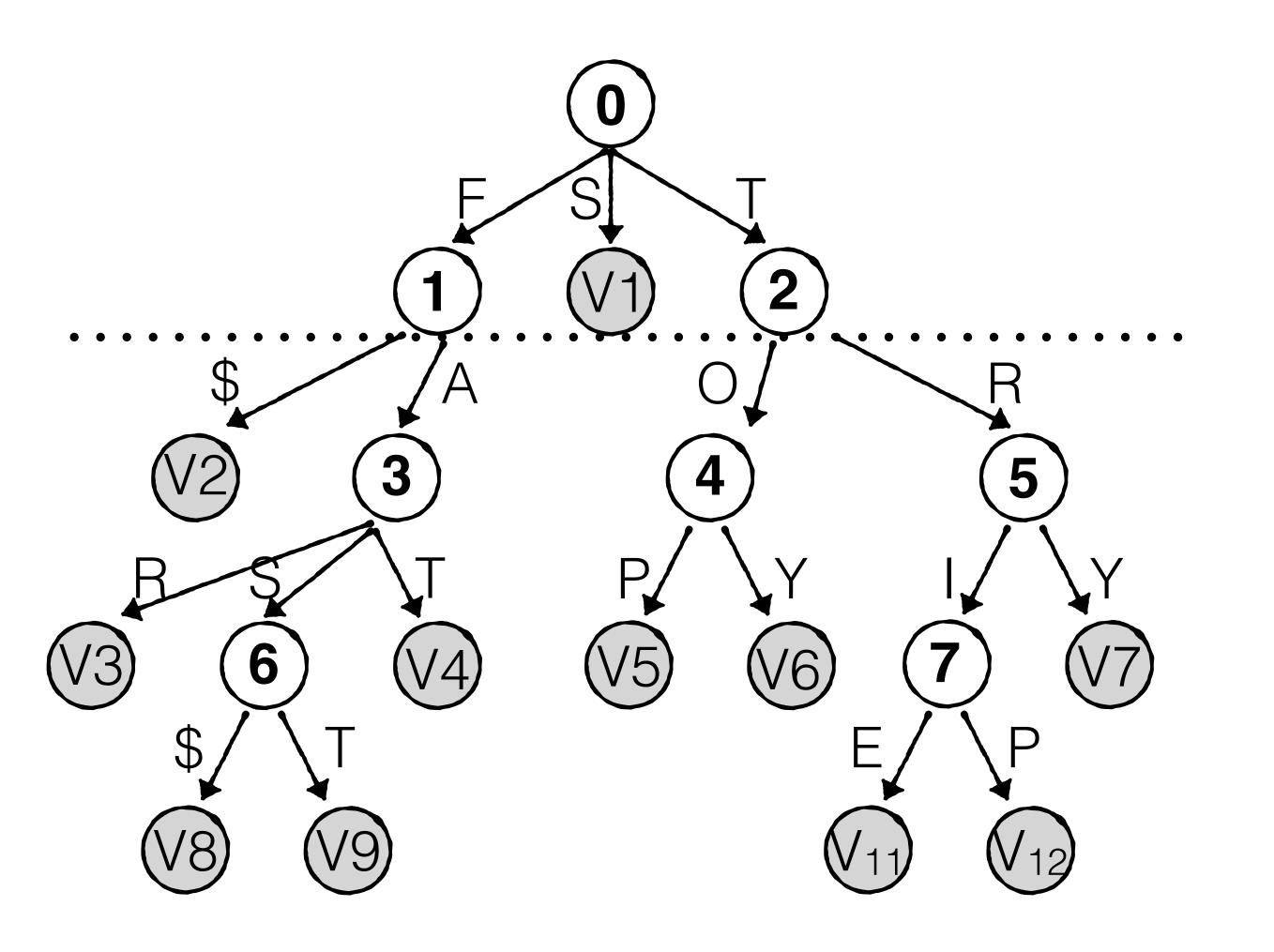


### A prefix is always represented by an internal node



# Label internal nodes in breadth-first order for clarity



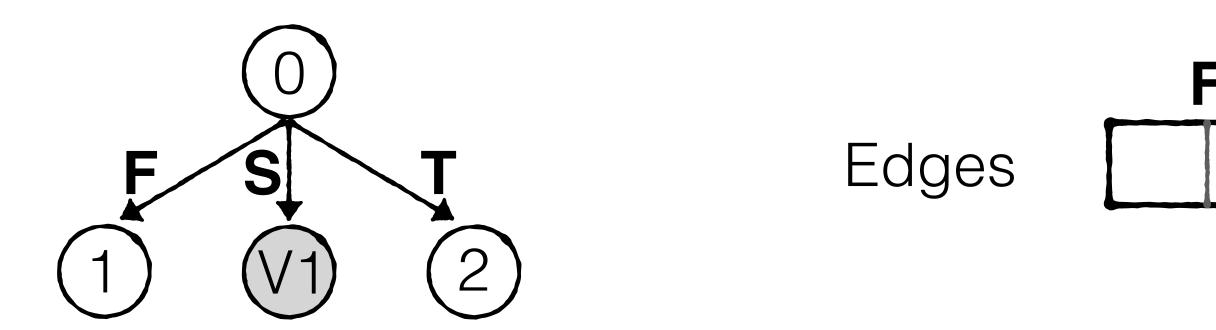


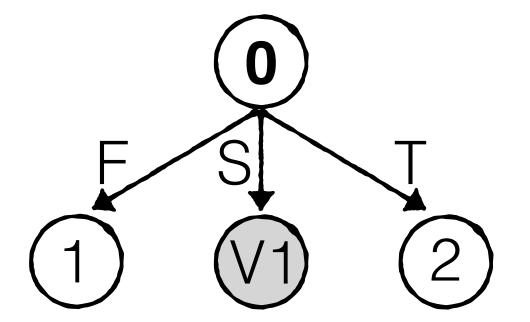
## Dense encoding

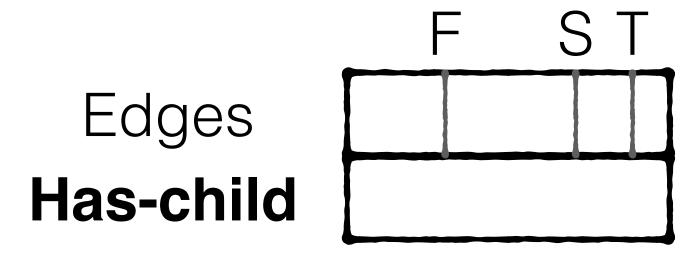
Sparse encoding

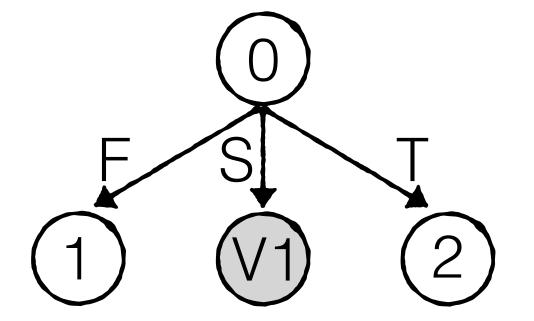


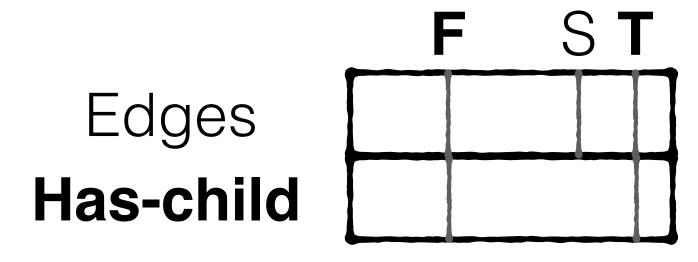
### One bit set for every connected character

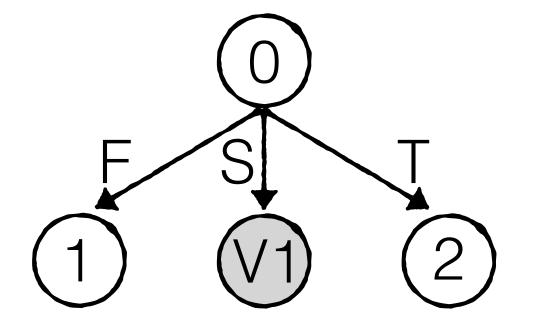


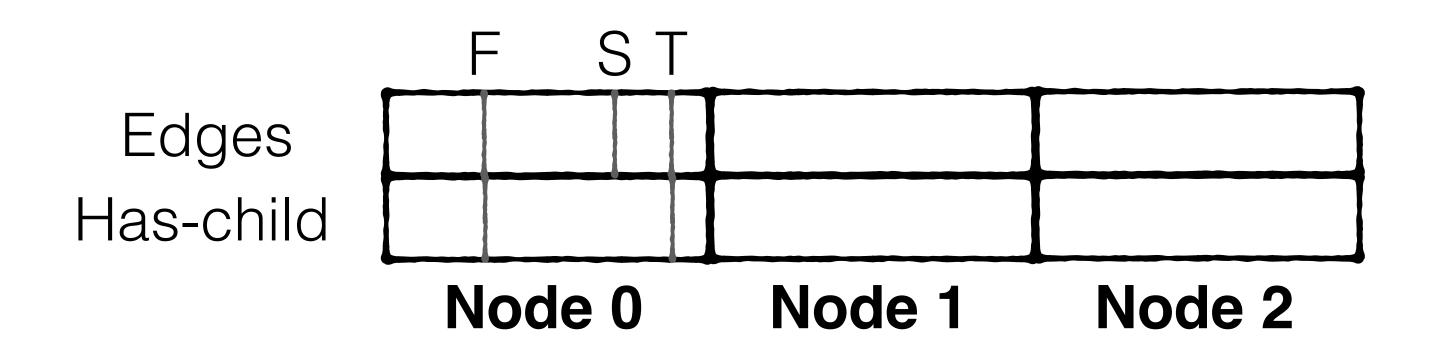




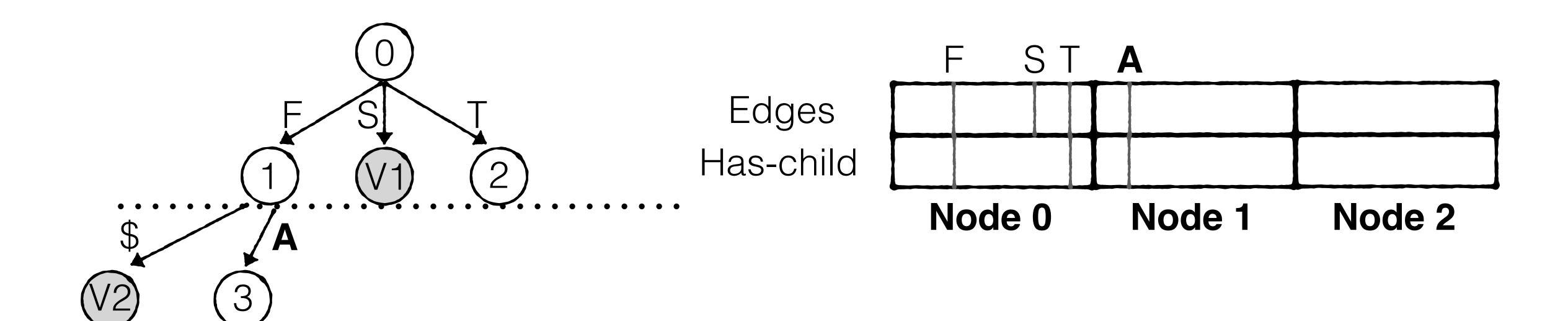




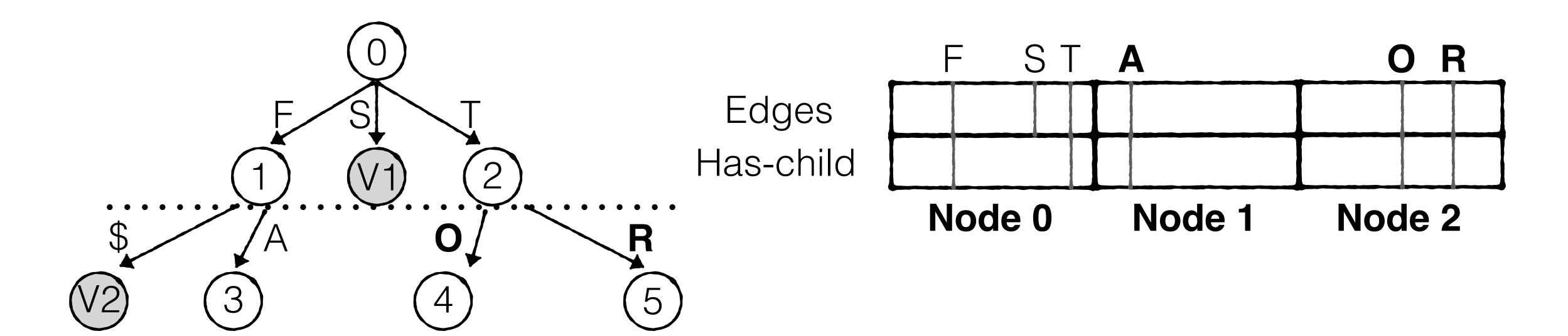




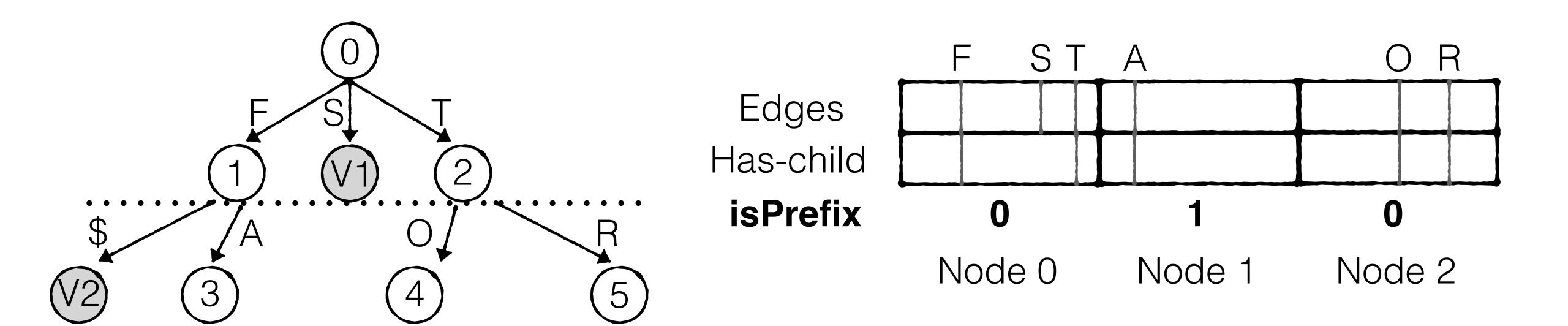
Continue in breadth-first order for each internal node

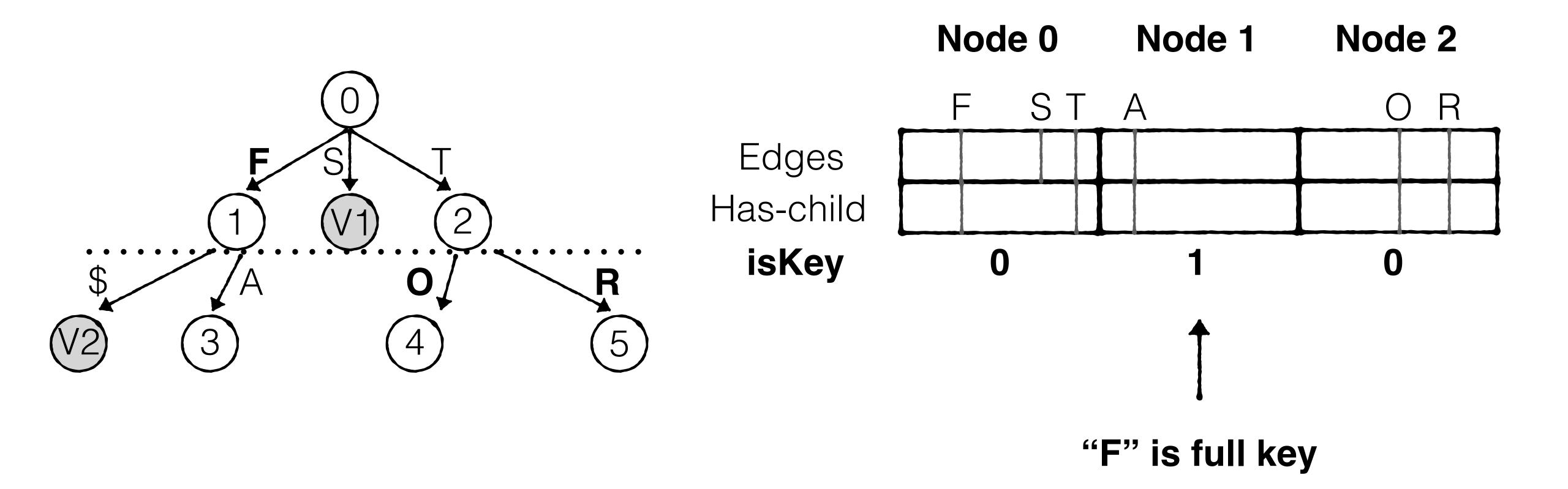


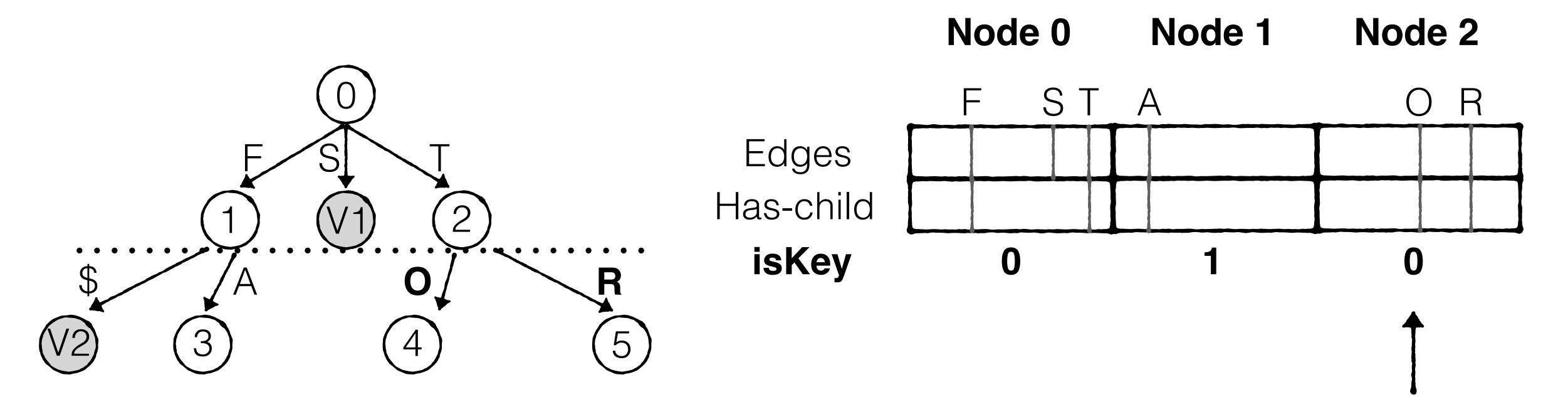
Continue in breadth-first order for each internal node



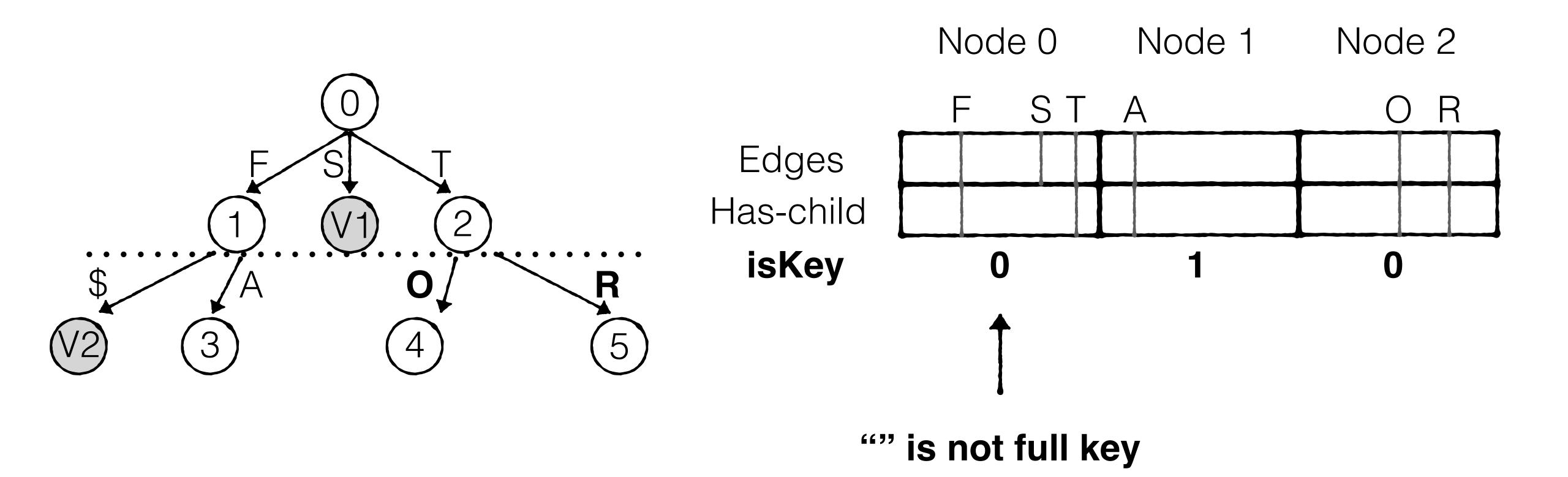
Continue in breadth-first order for each internal node

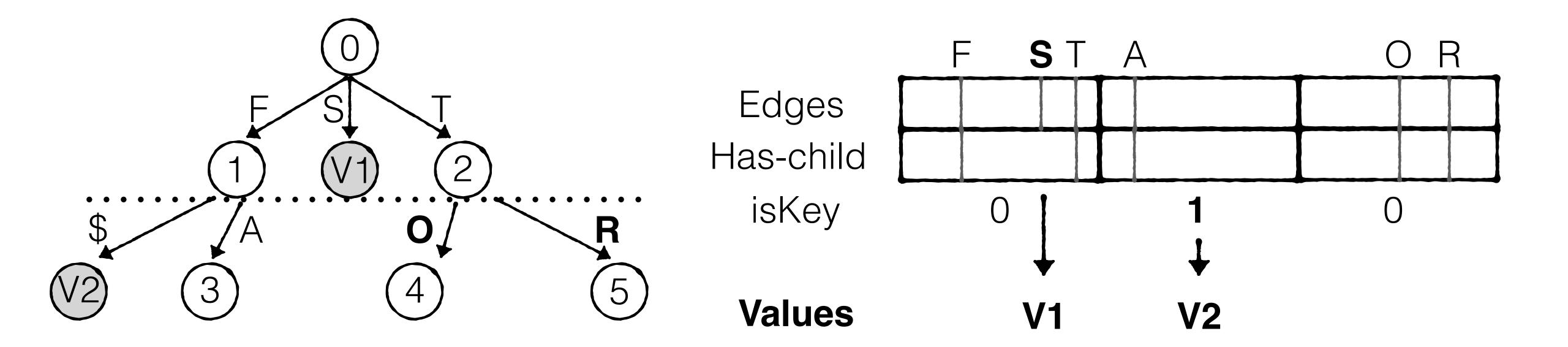


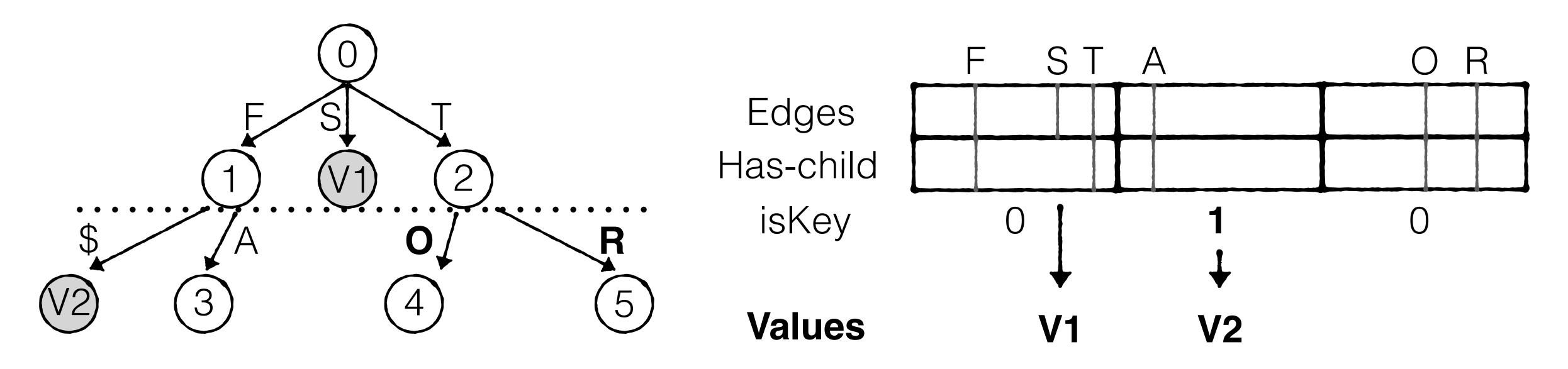




"T" is not full key

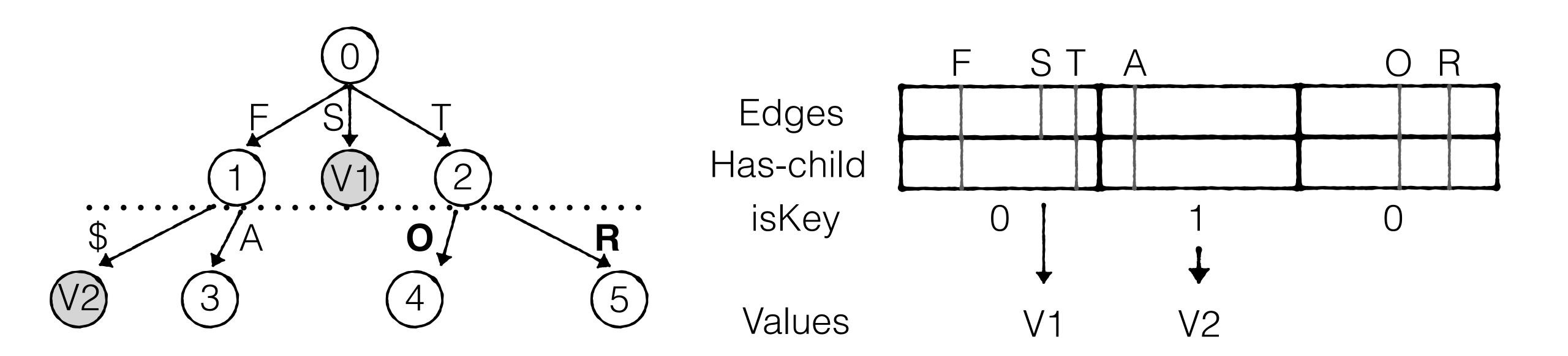


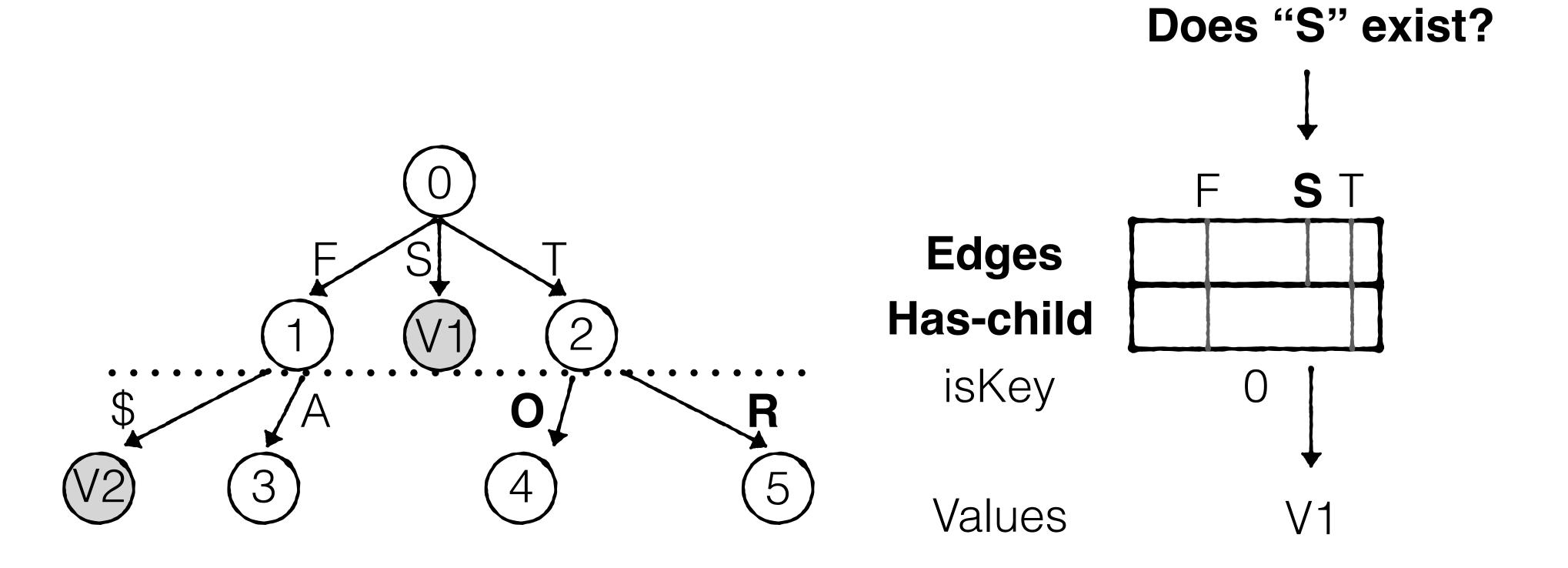


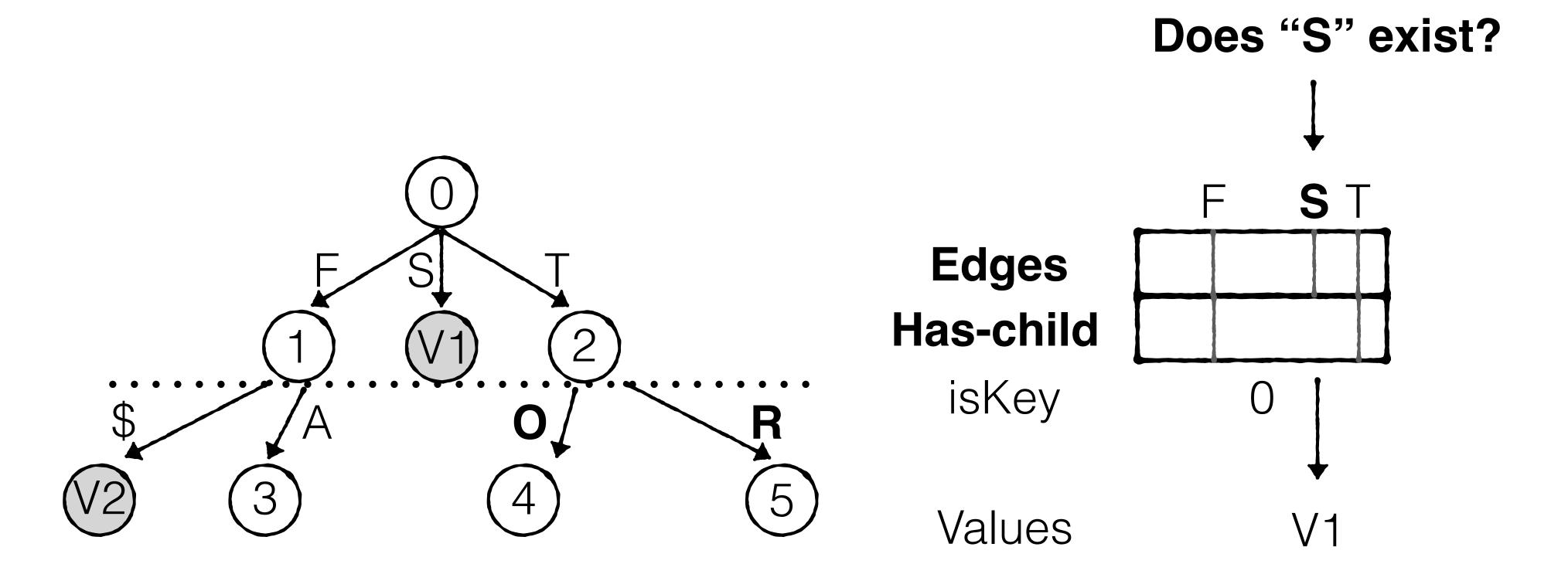


Stored contiguously in an array

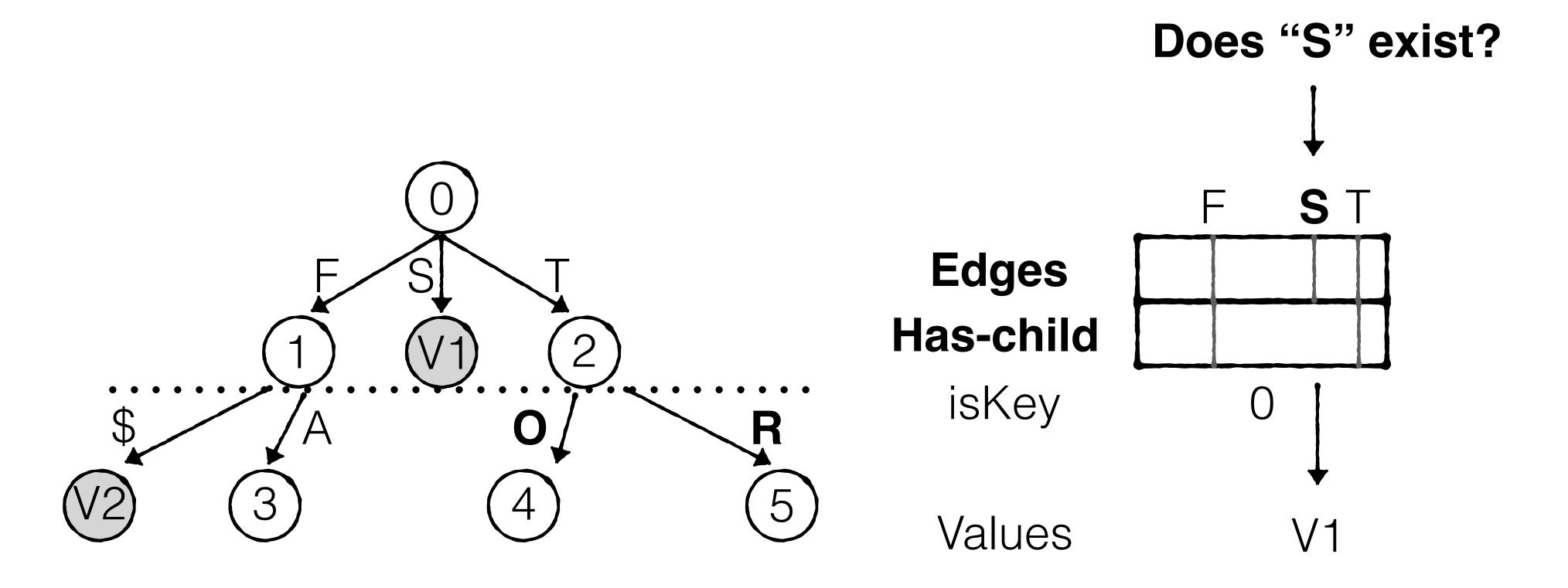
Does "S" exist?





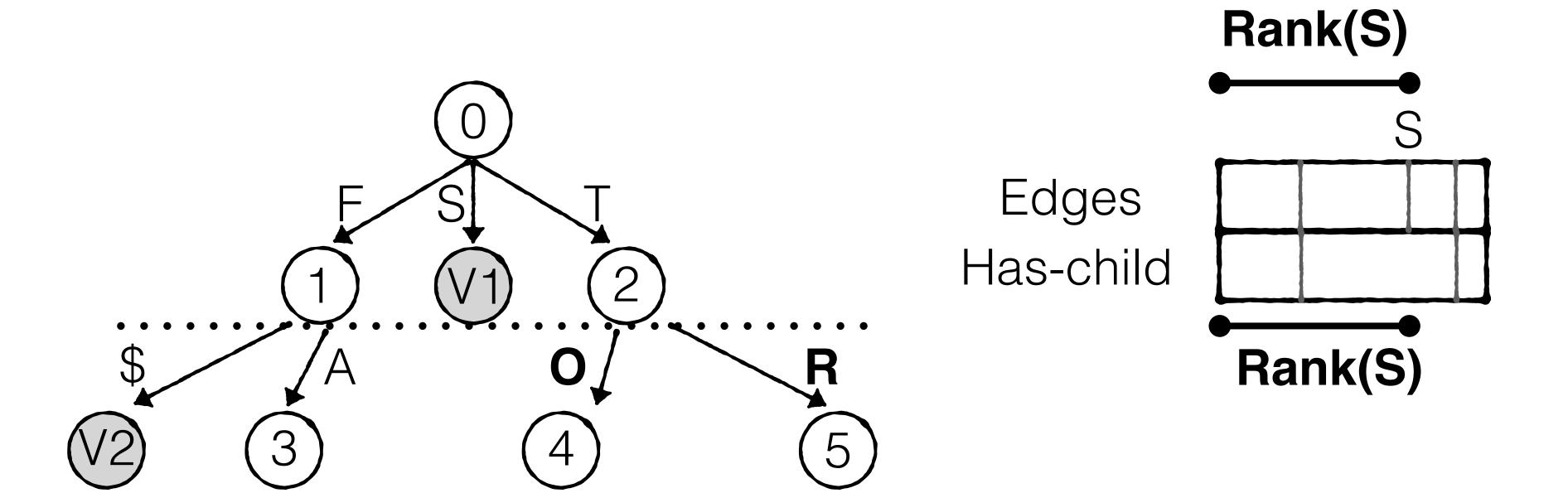


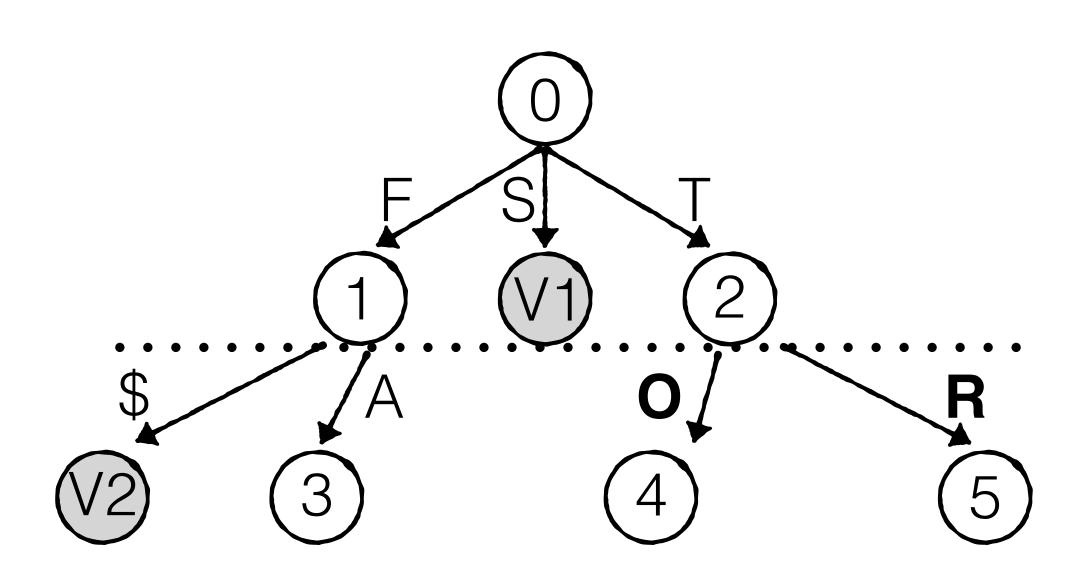
We observe that S exists and has no children, meaning it must be a full key



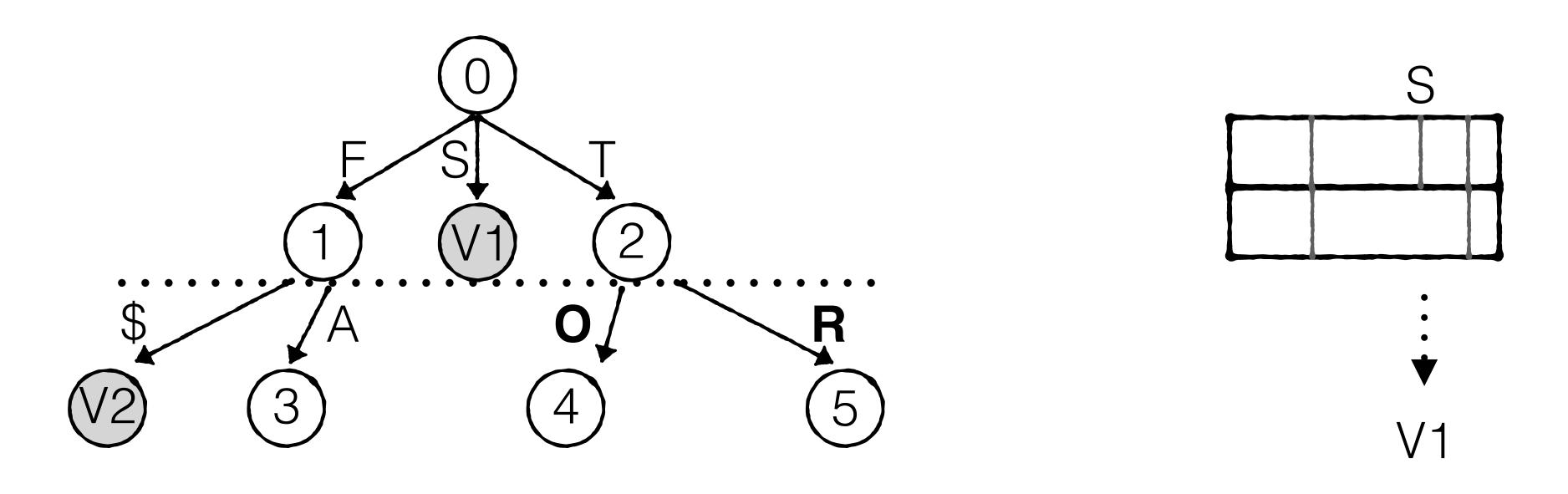
We observe that S exists and has no children, meaning it must be a full key

## How to locate its value?

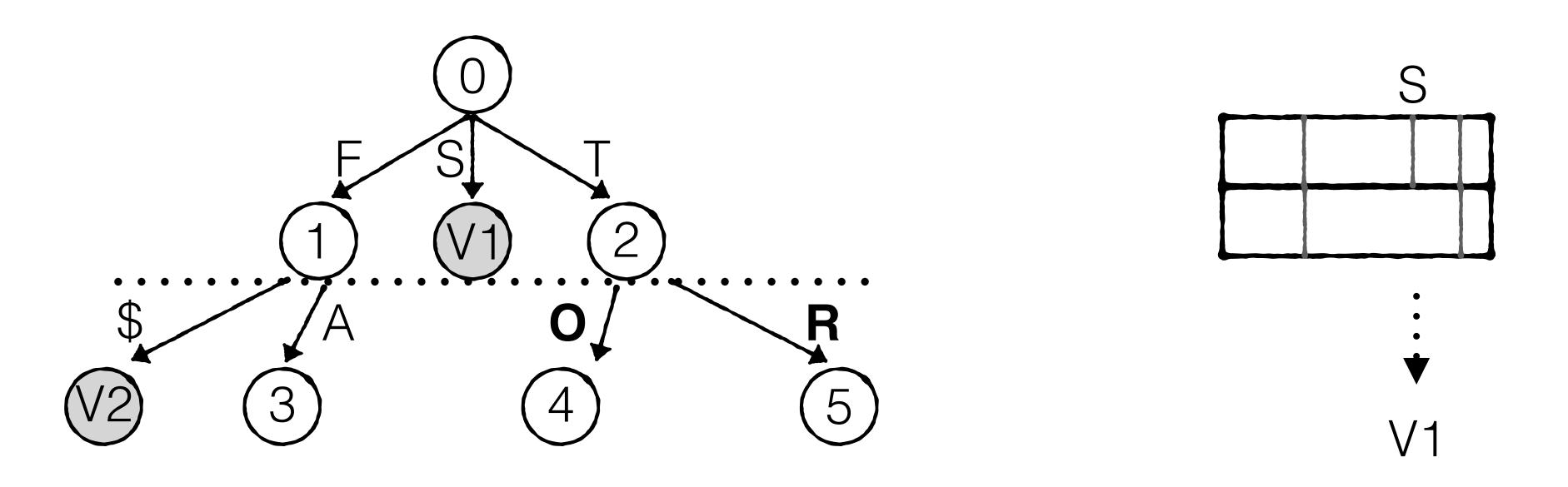




S value offset = Rank(Edges, S) - Rank(Has-child, S) = 0

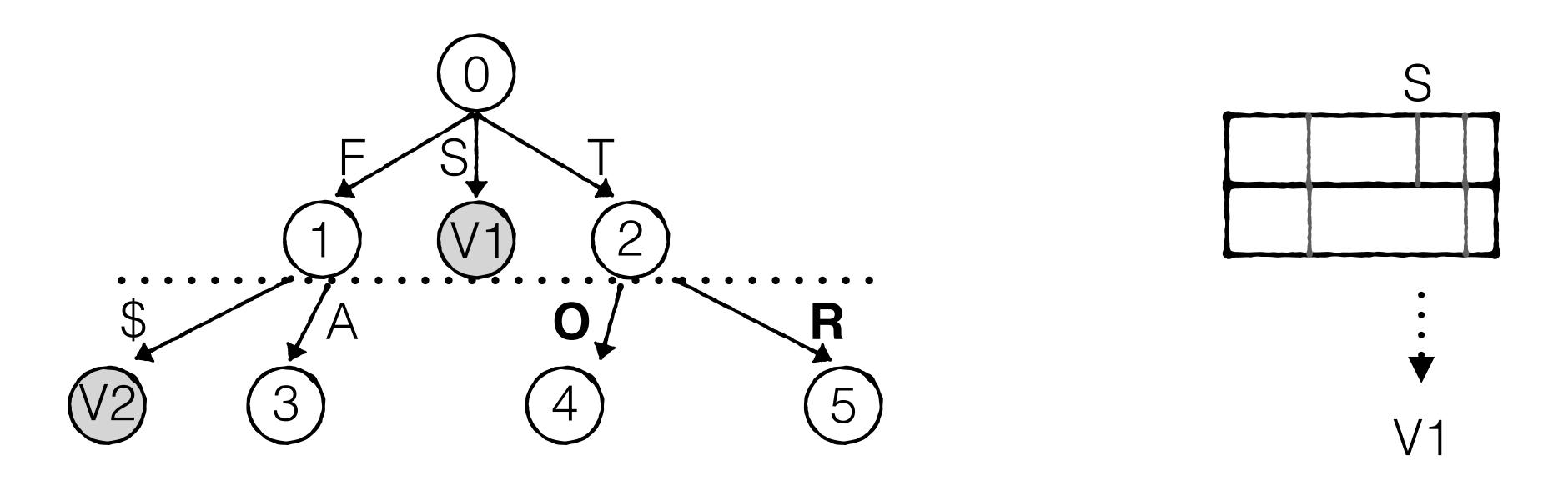


S value offset = Rank(Edges, S) - Rank(Has-child, S) = 0



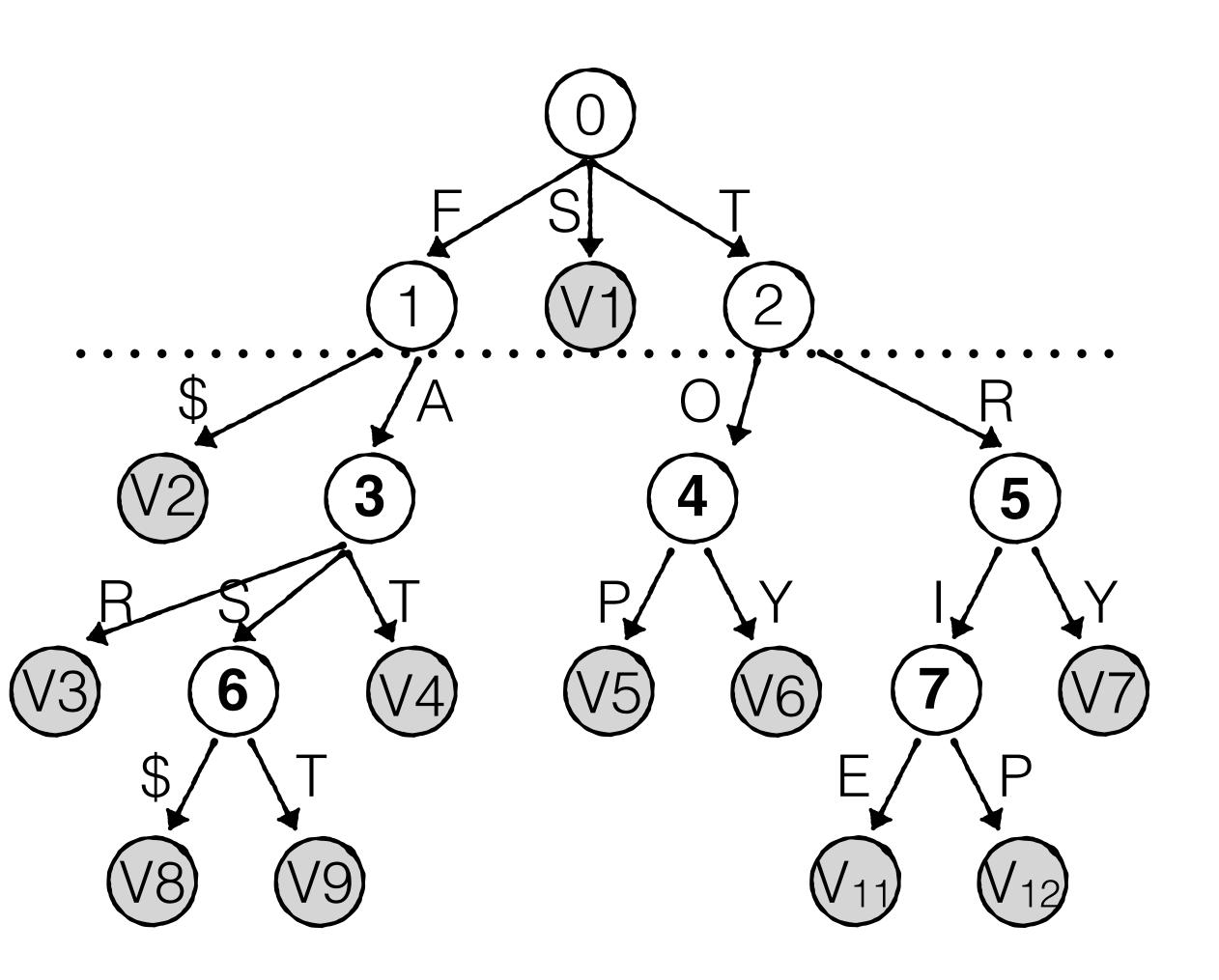
Constant time per each node access:)

S value offset = Rank(Edges, S) - Rank(Has-child, S) = 0



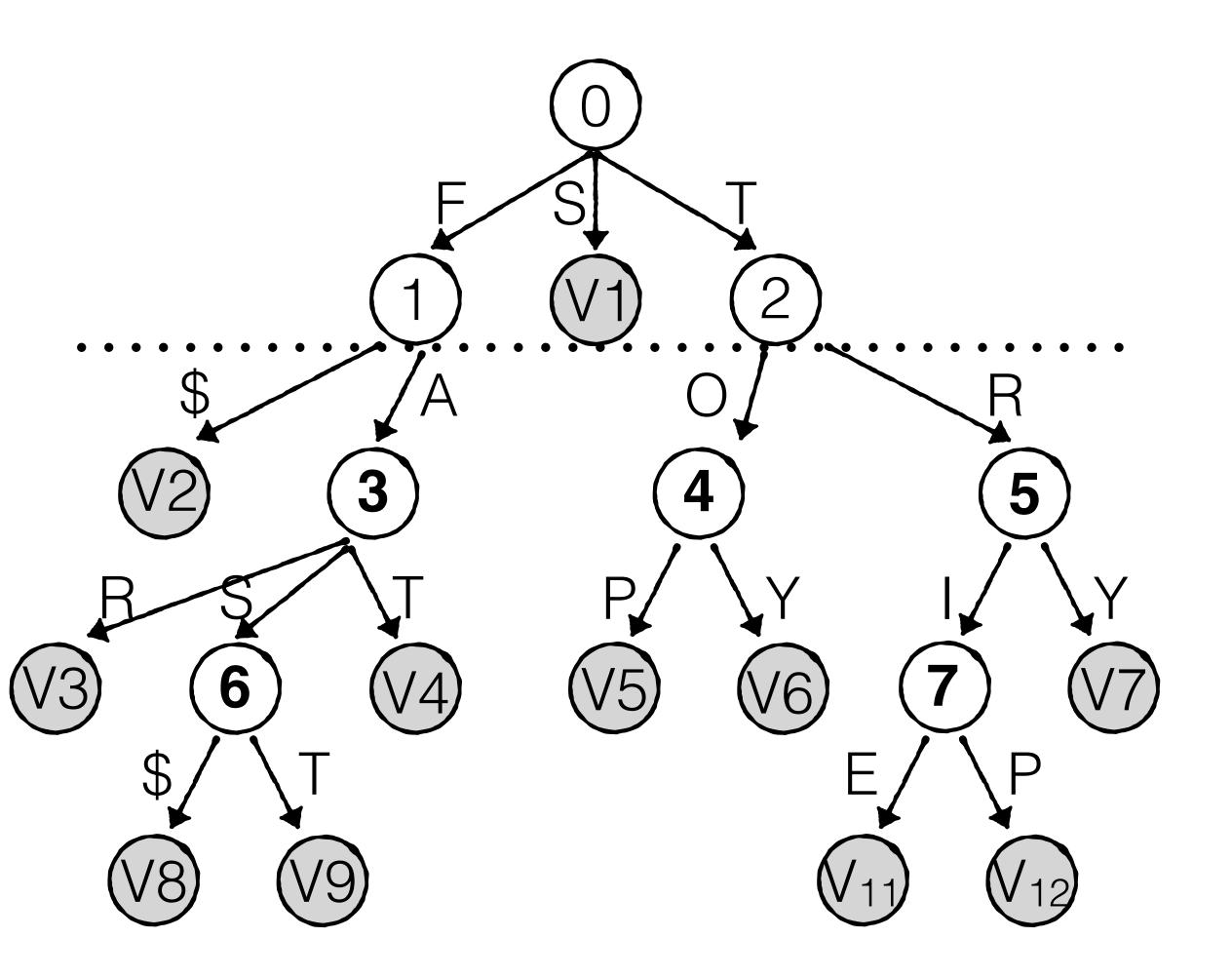
Constant time per each node access:)

Few bits per entry when base nodes have many edges

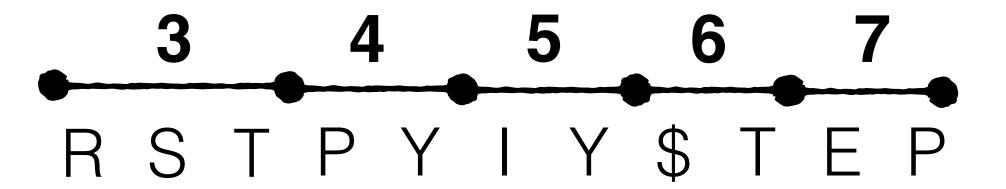


## Internal Nodes in breadth first order

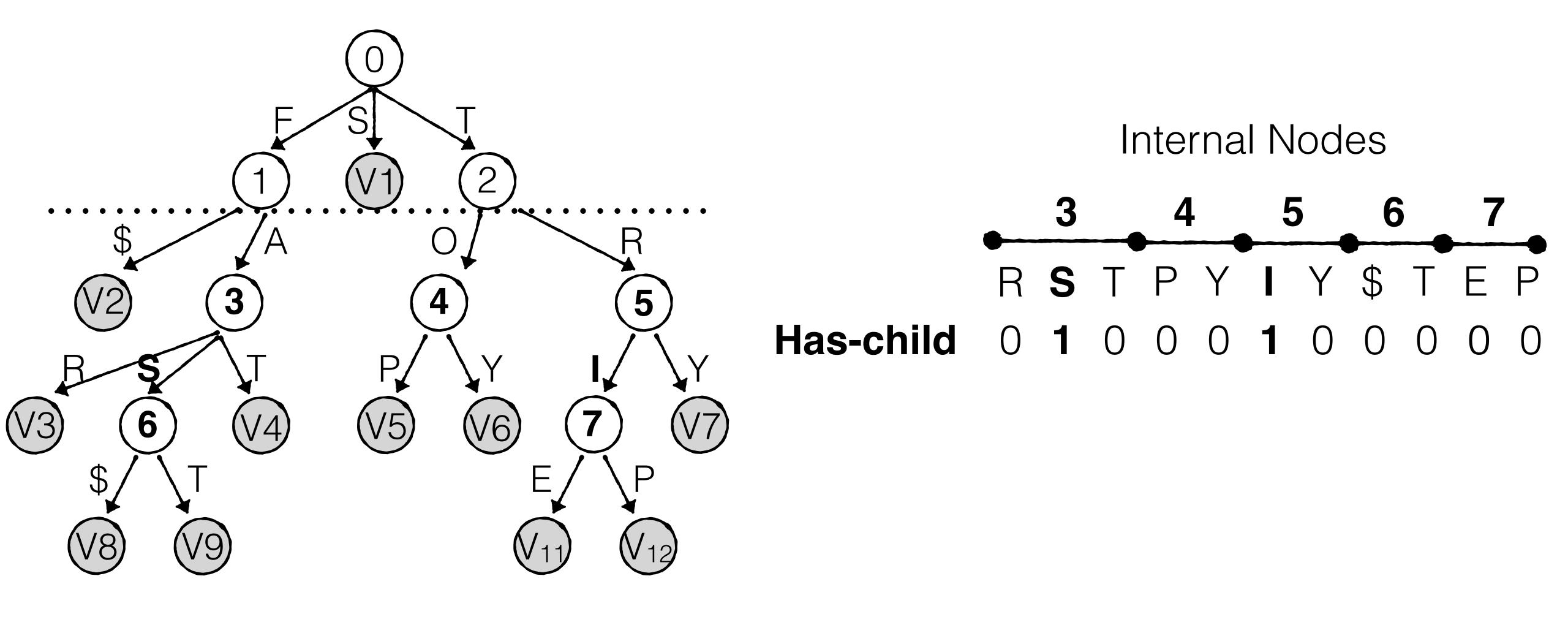
3 4 5 6 7

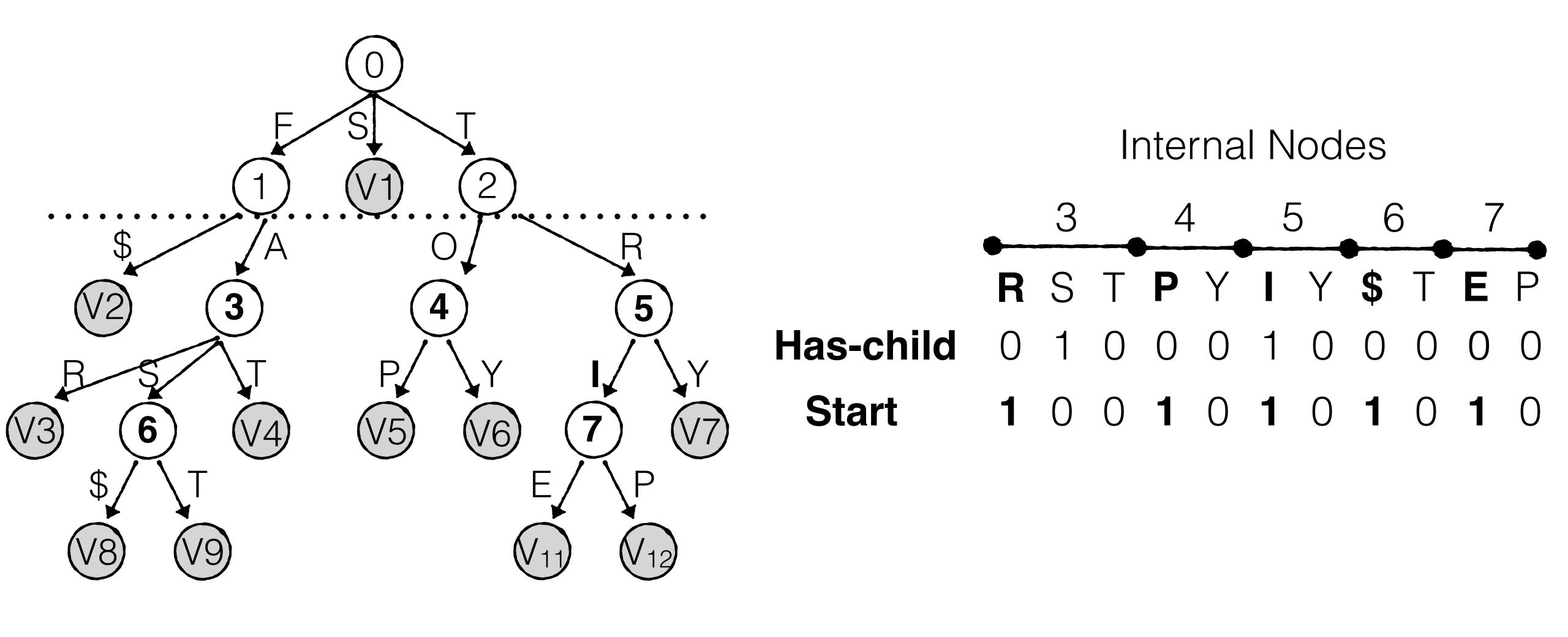


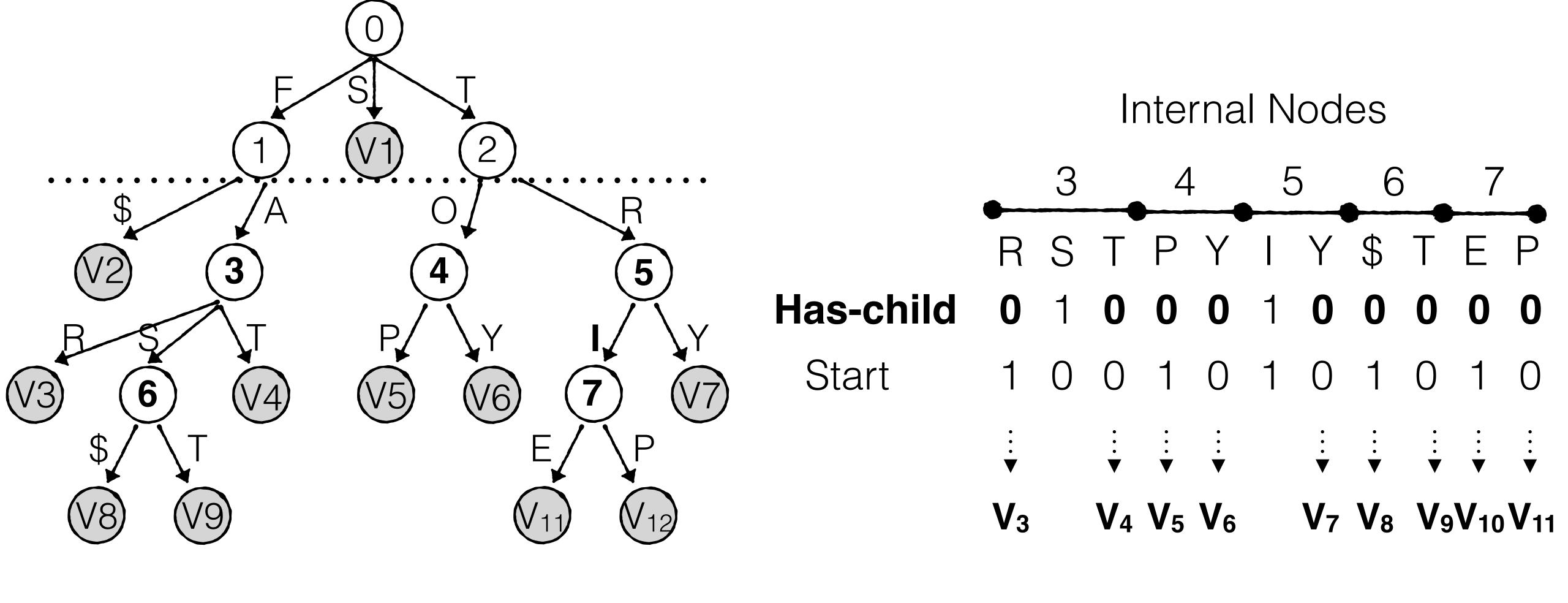
Internal Nodes in breadth first order

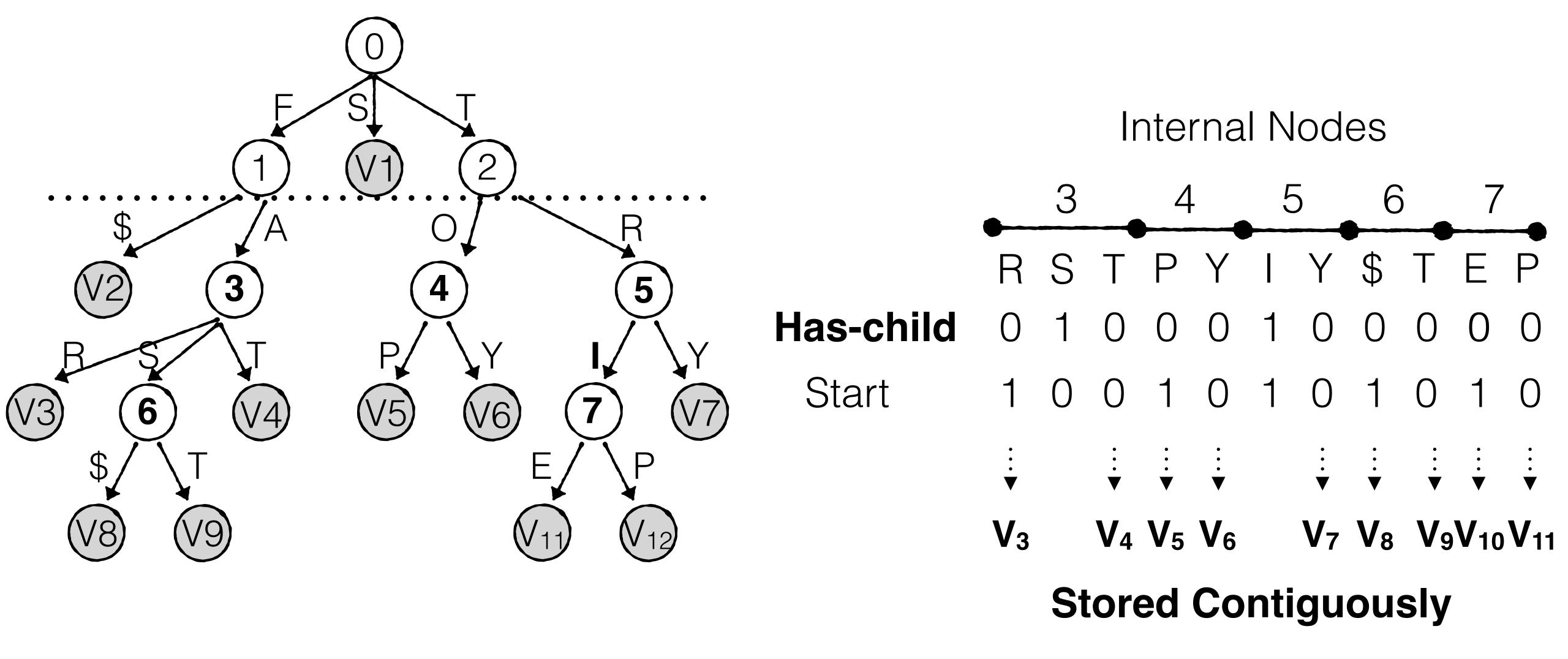


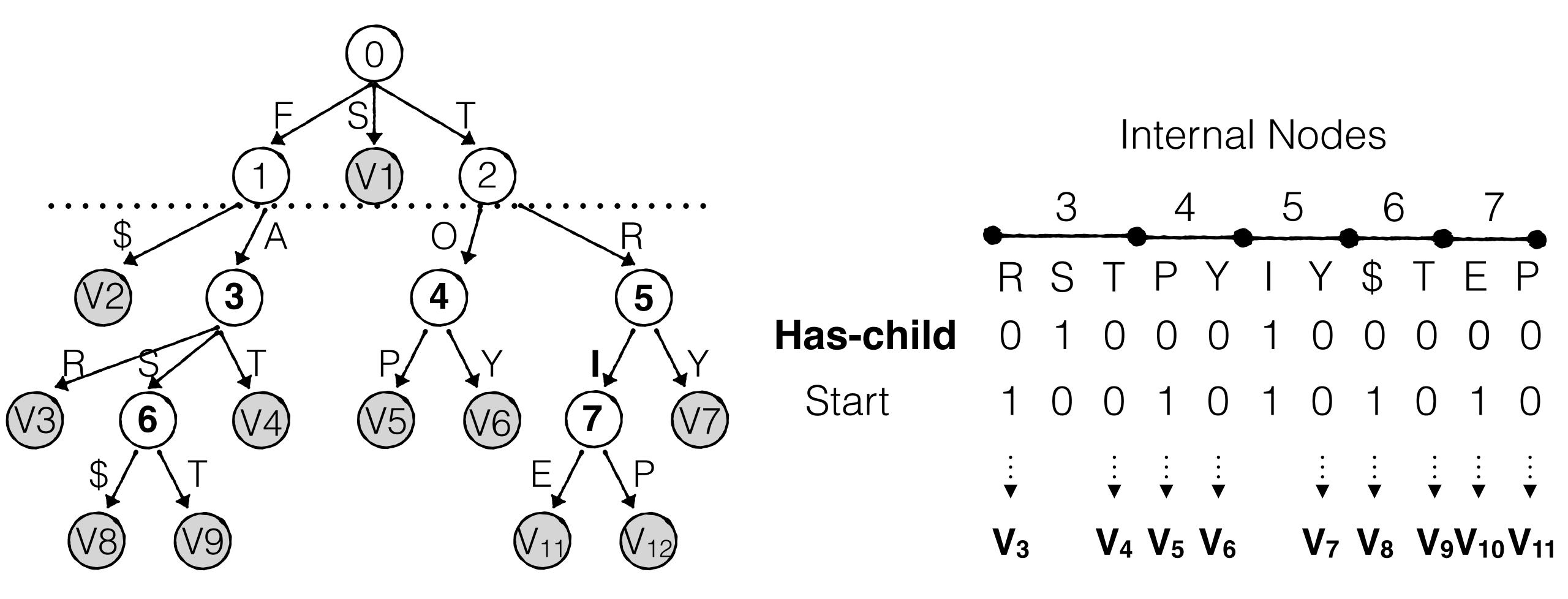
Store children contiguously



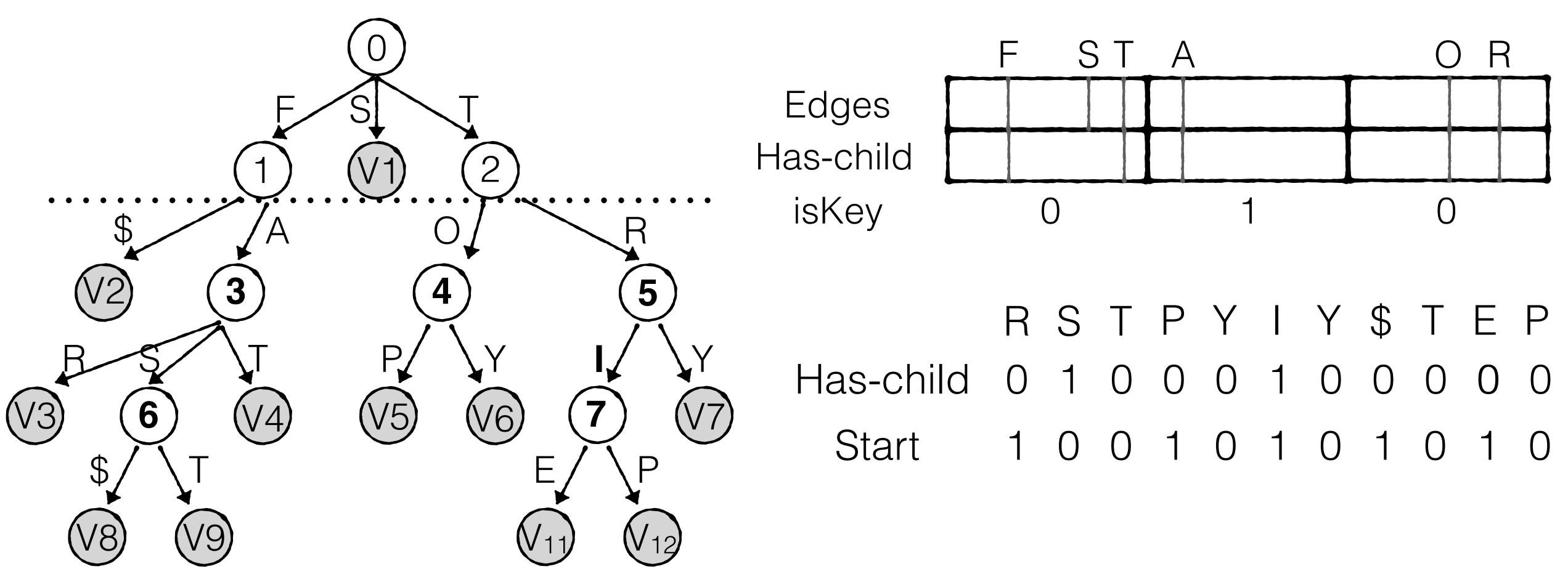




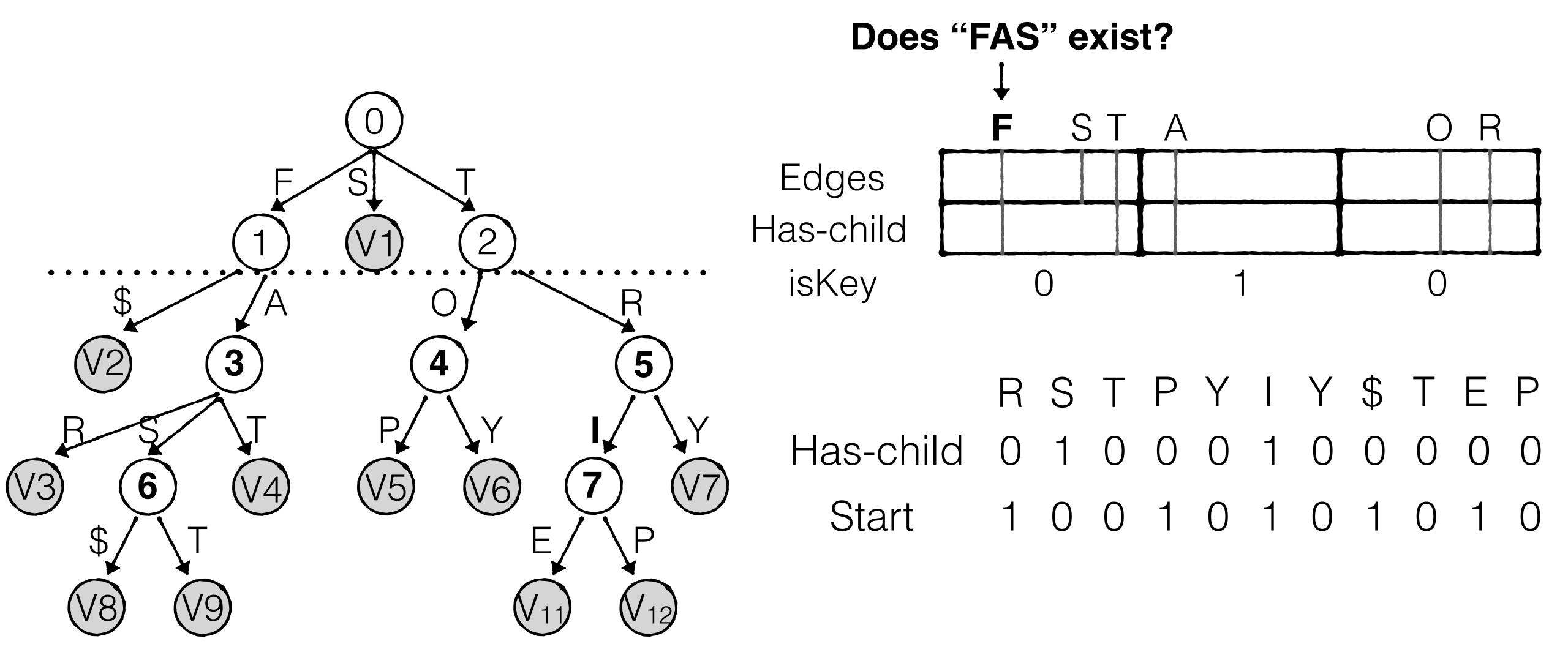


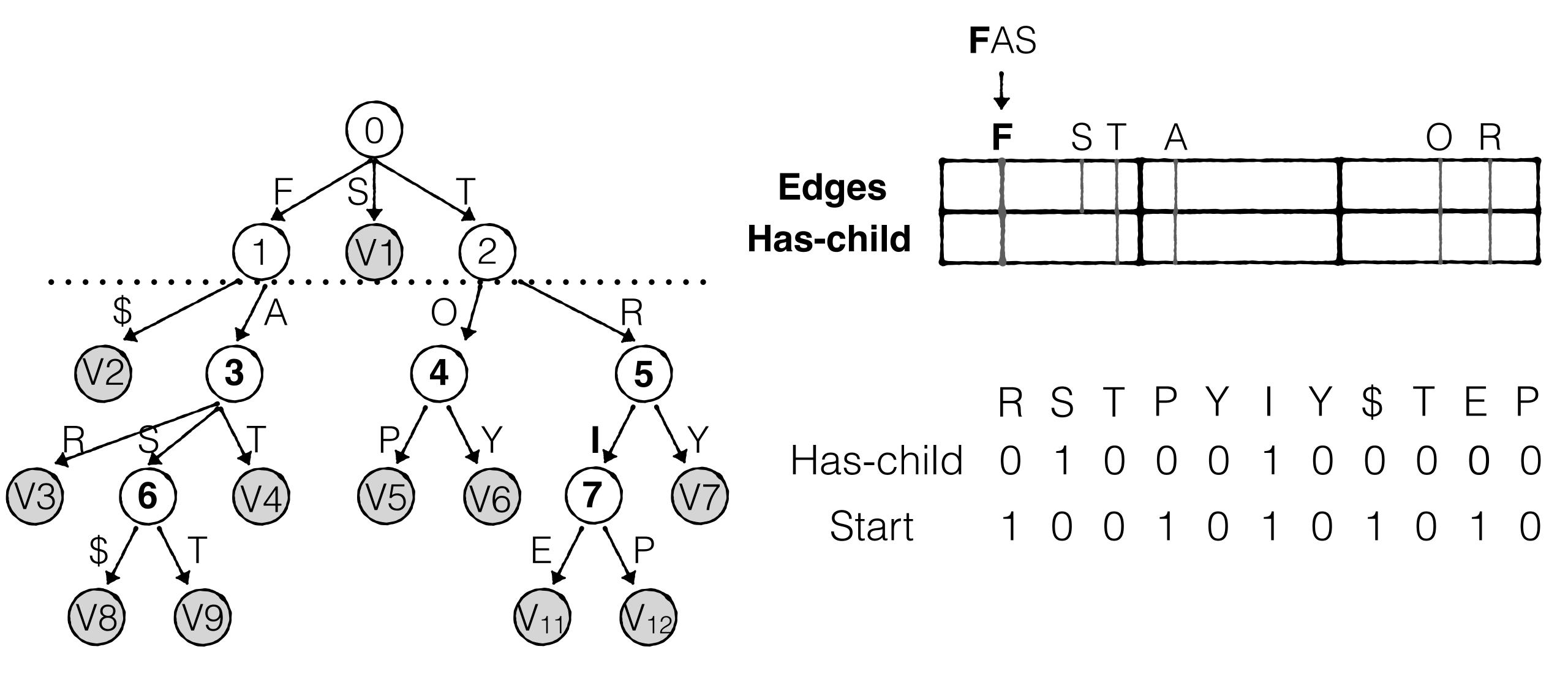


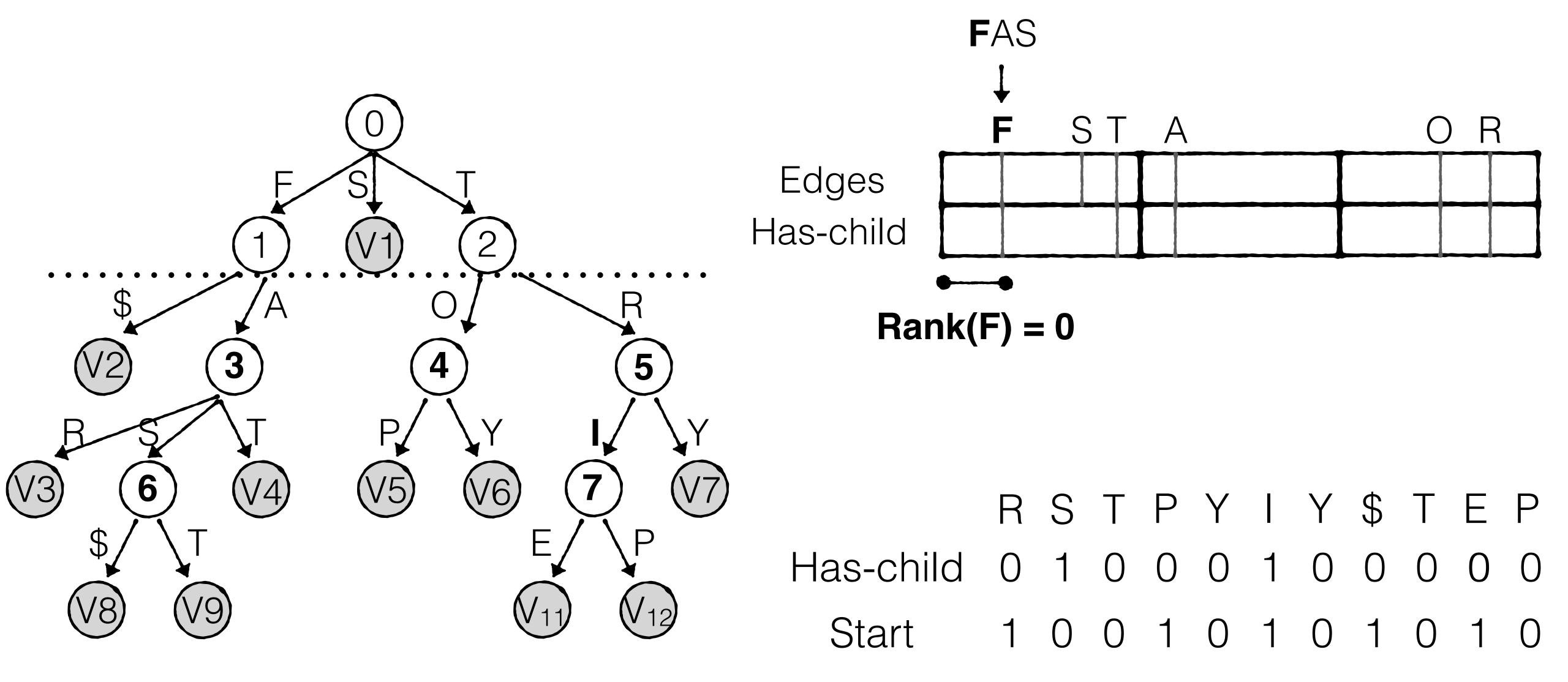
Each leaf connected to payload

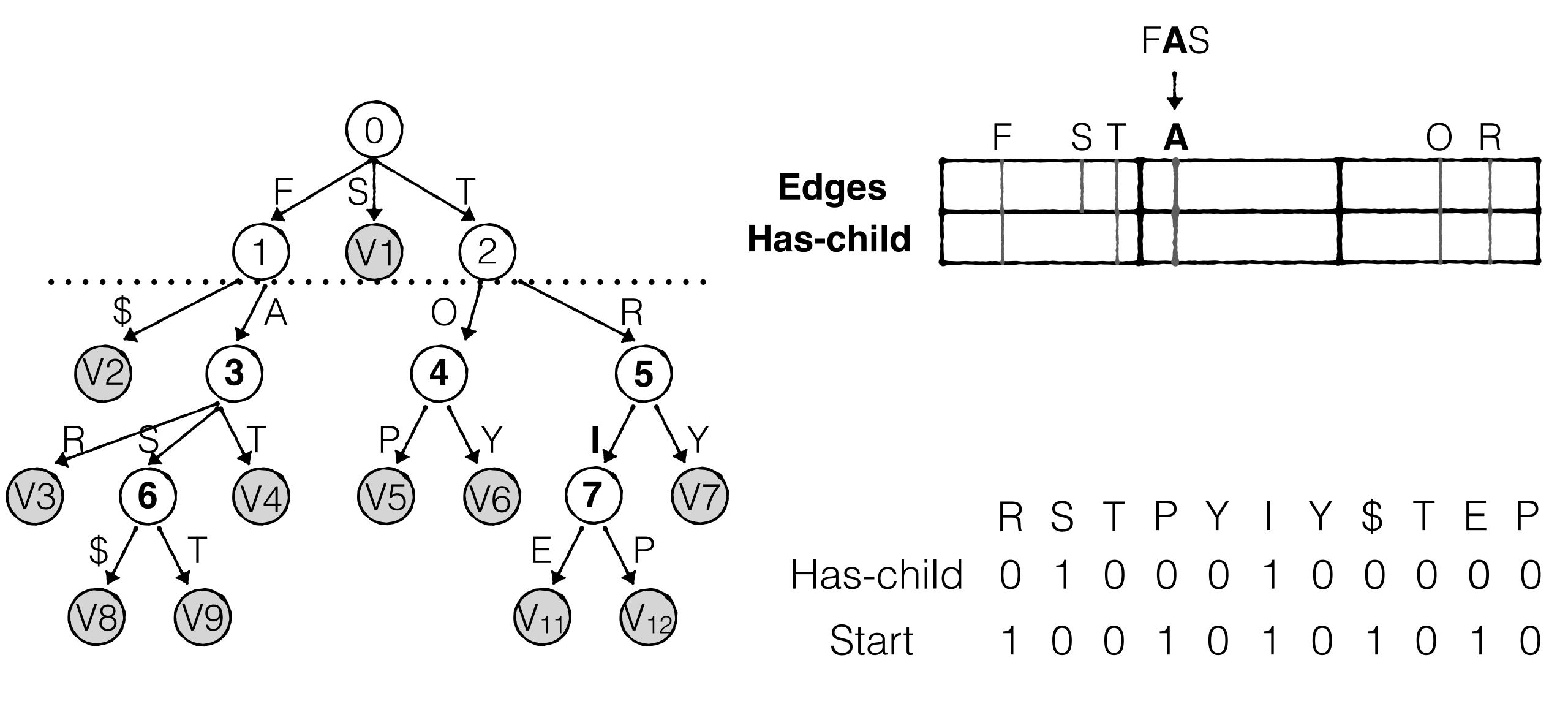


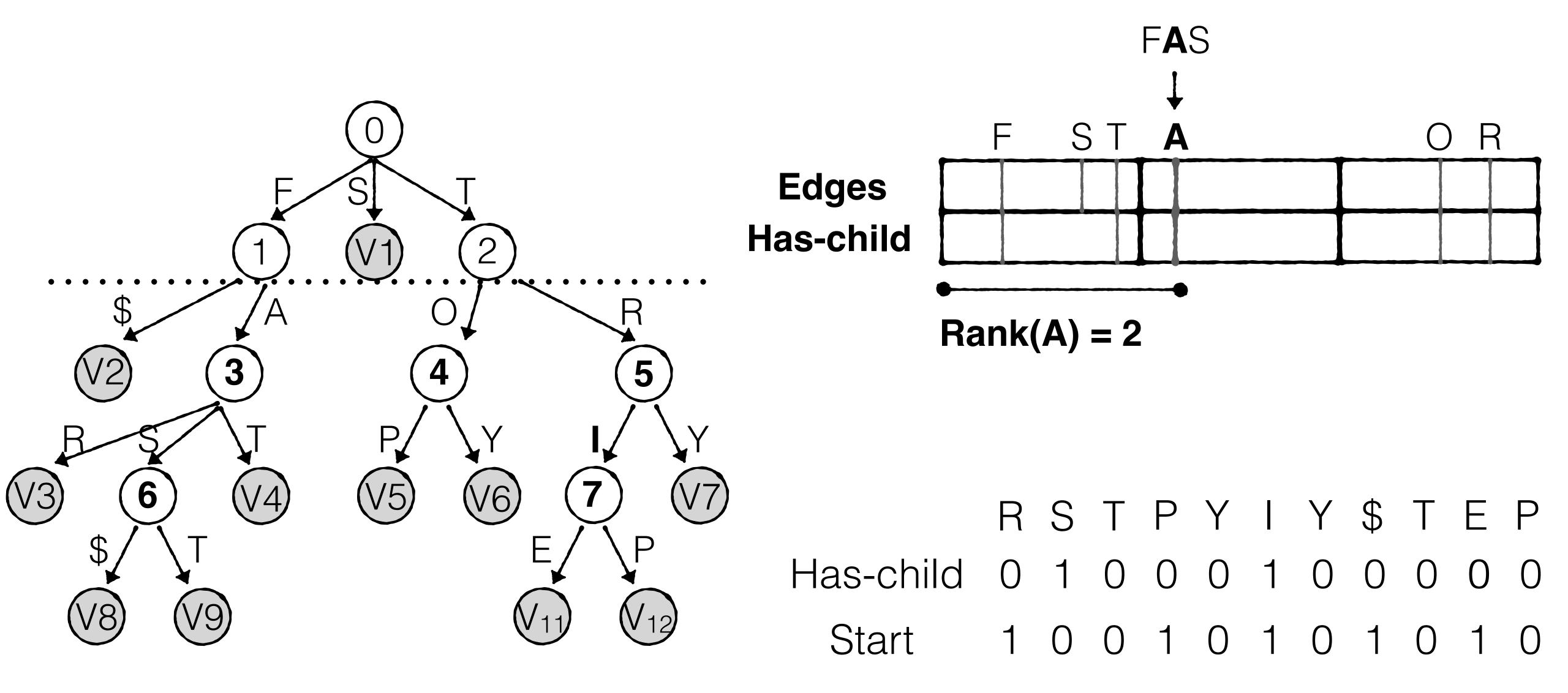
Whole filter:)

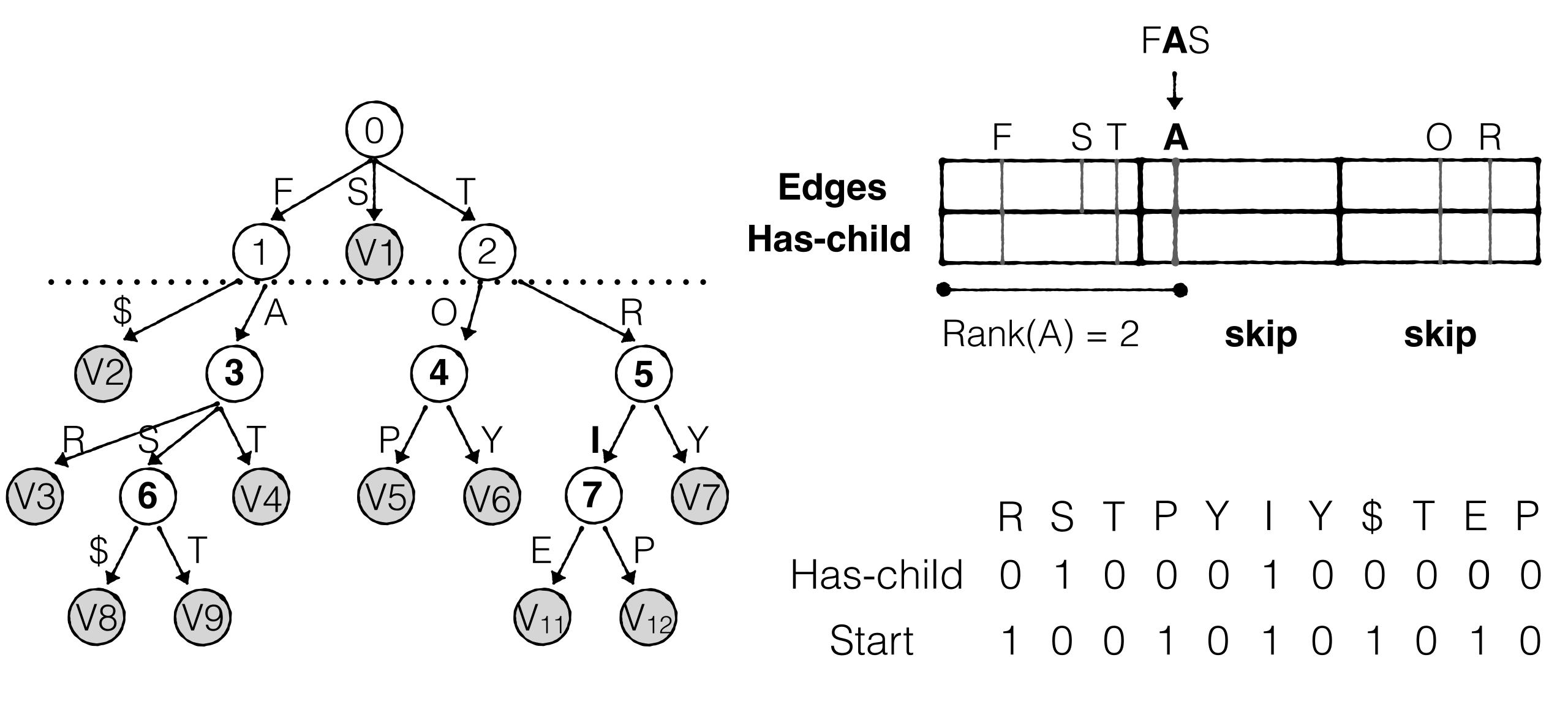


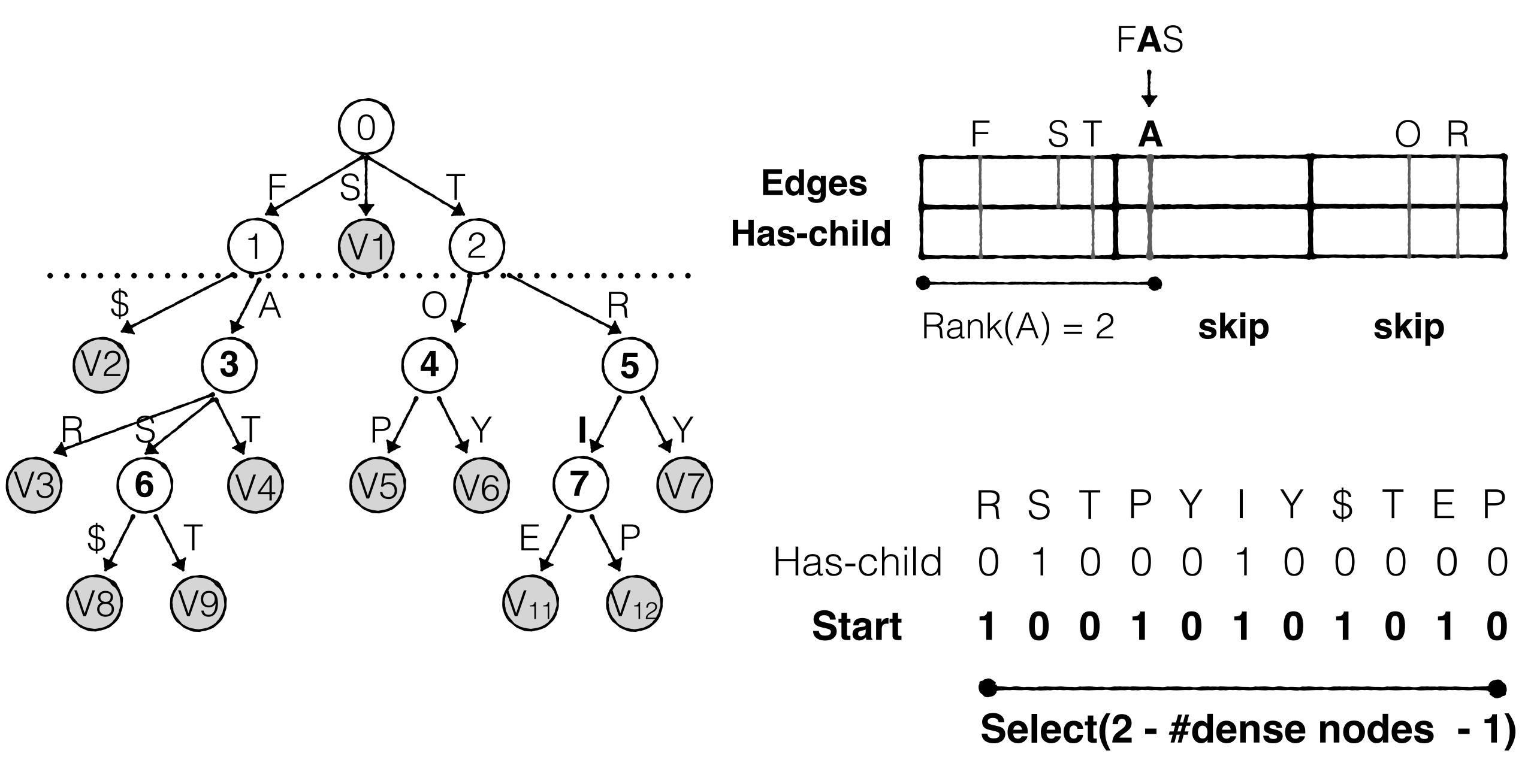


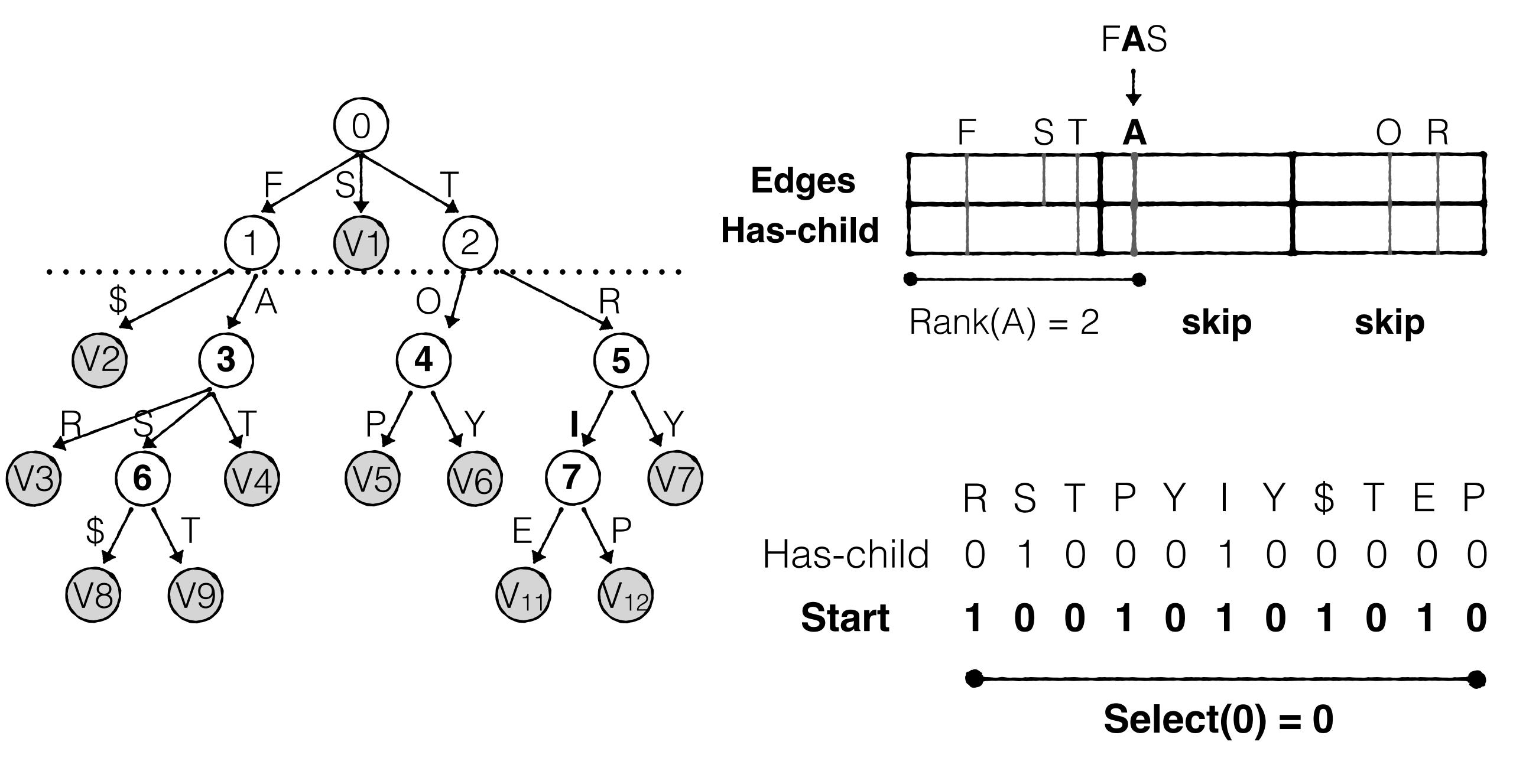


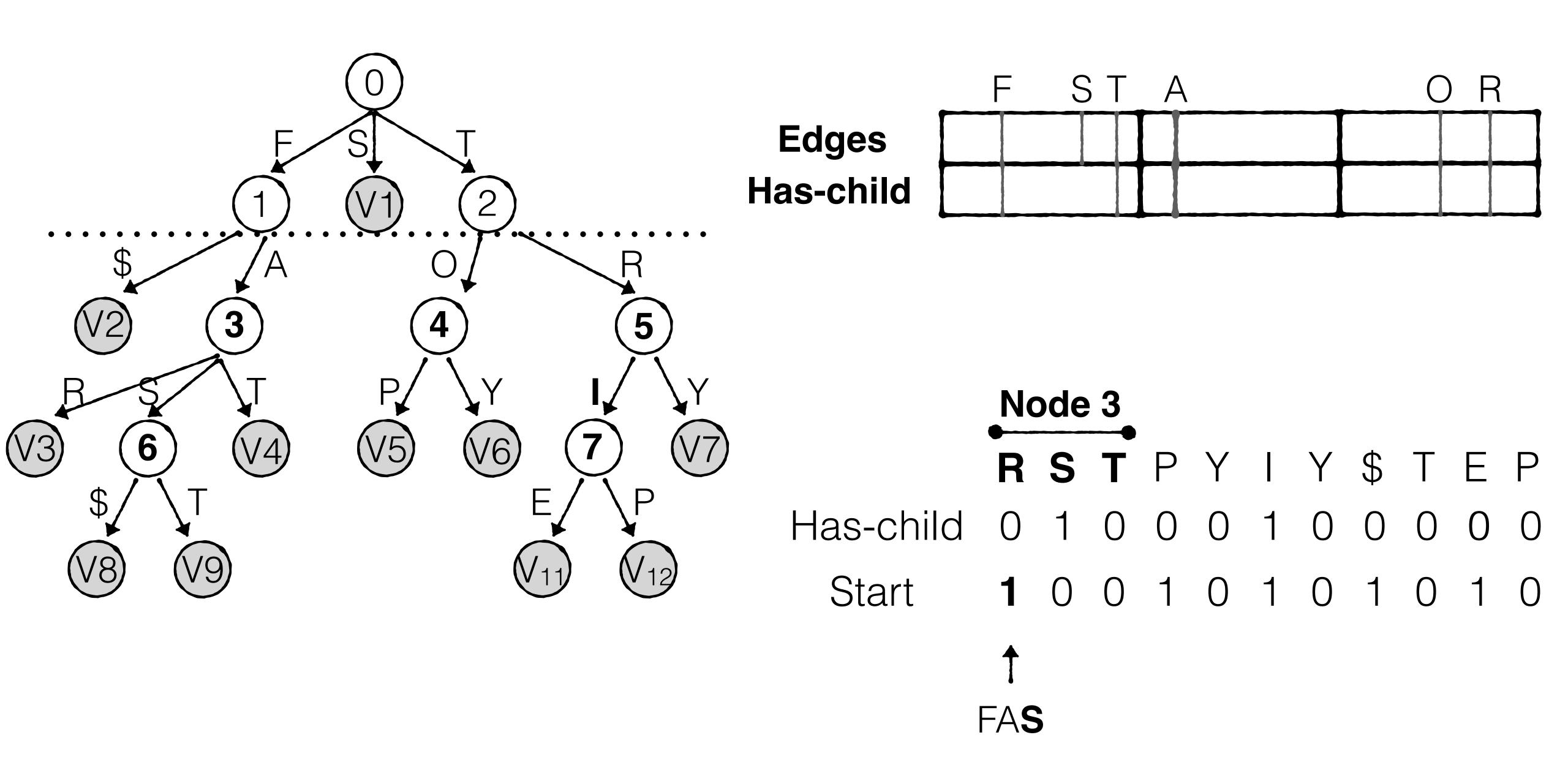


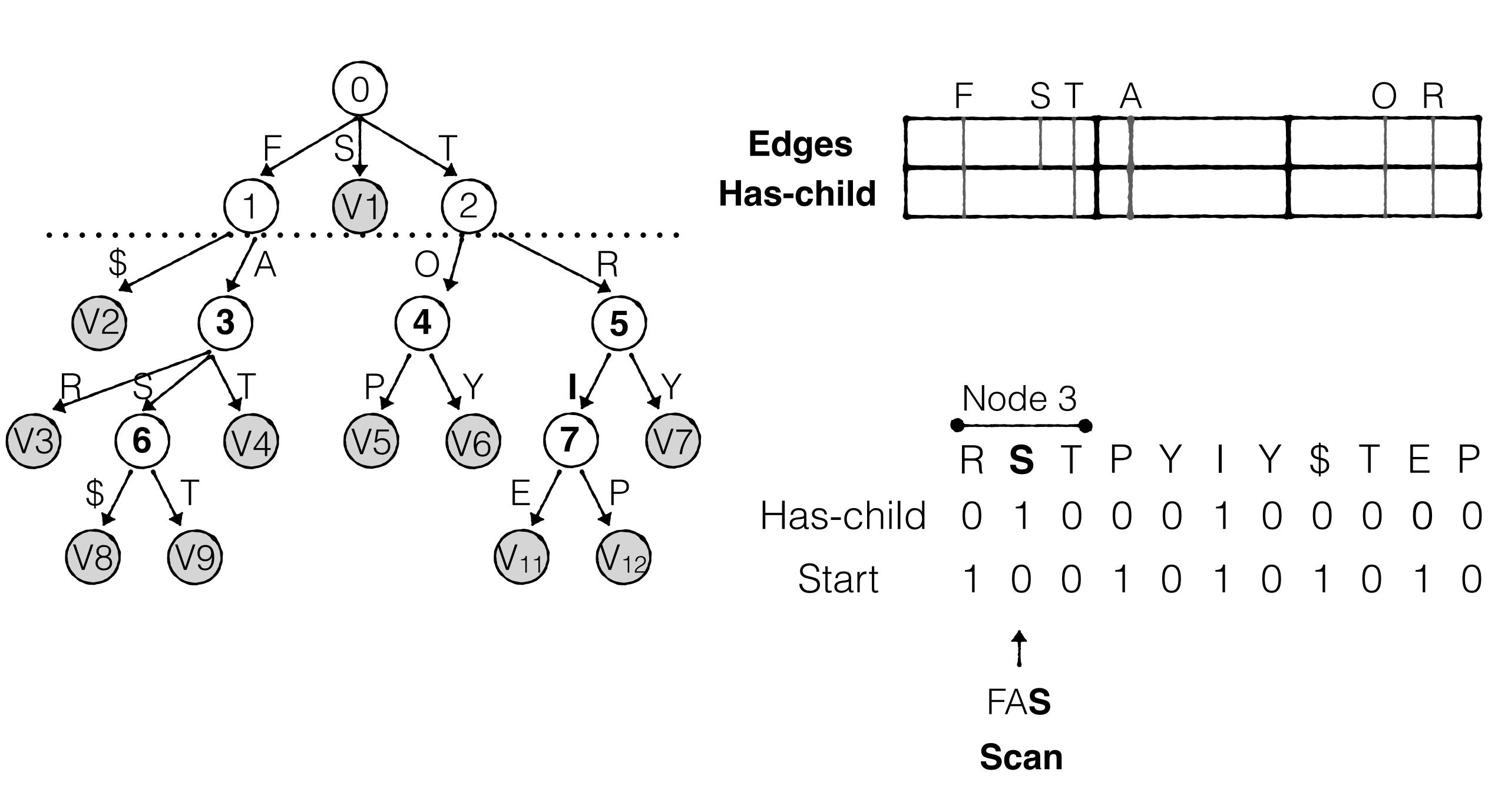


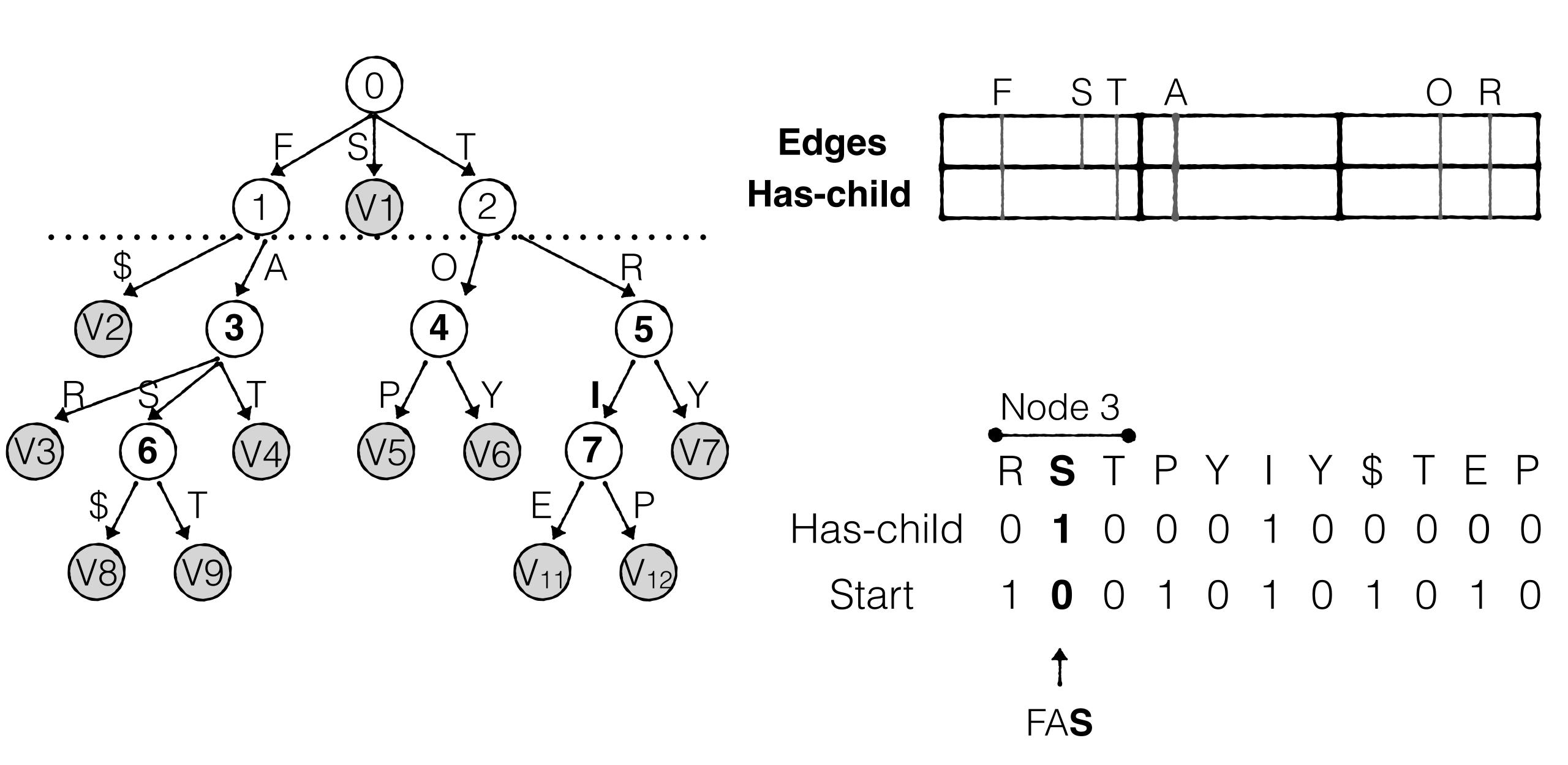


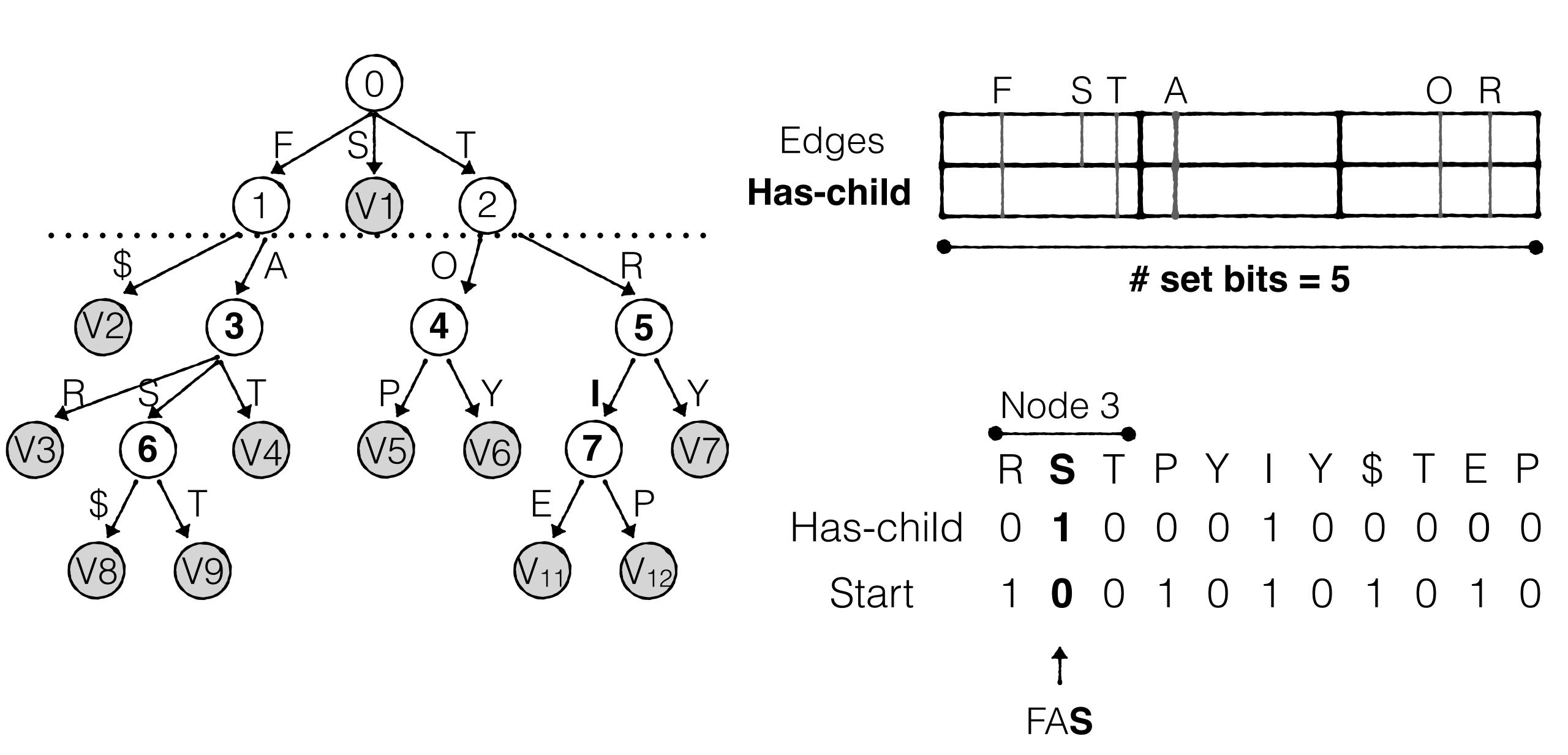


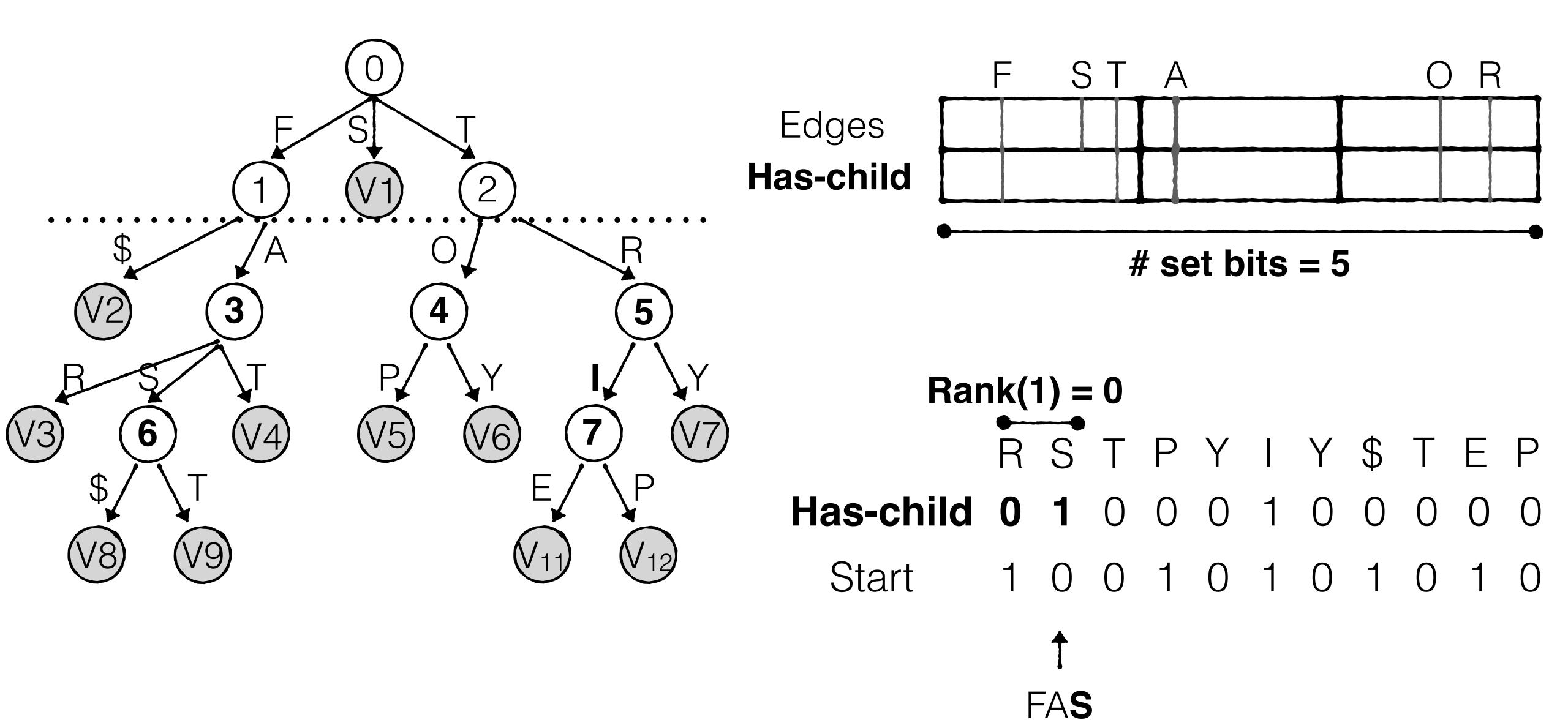


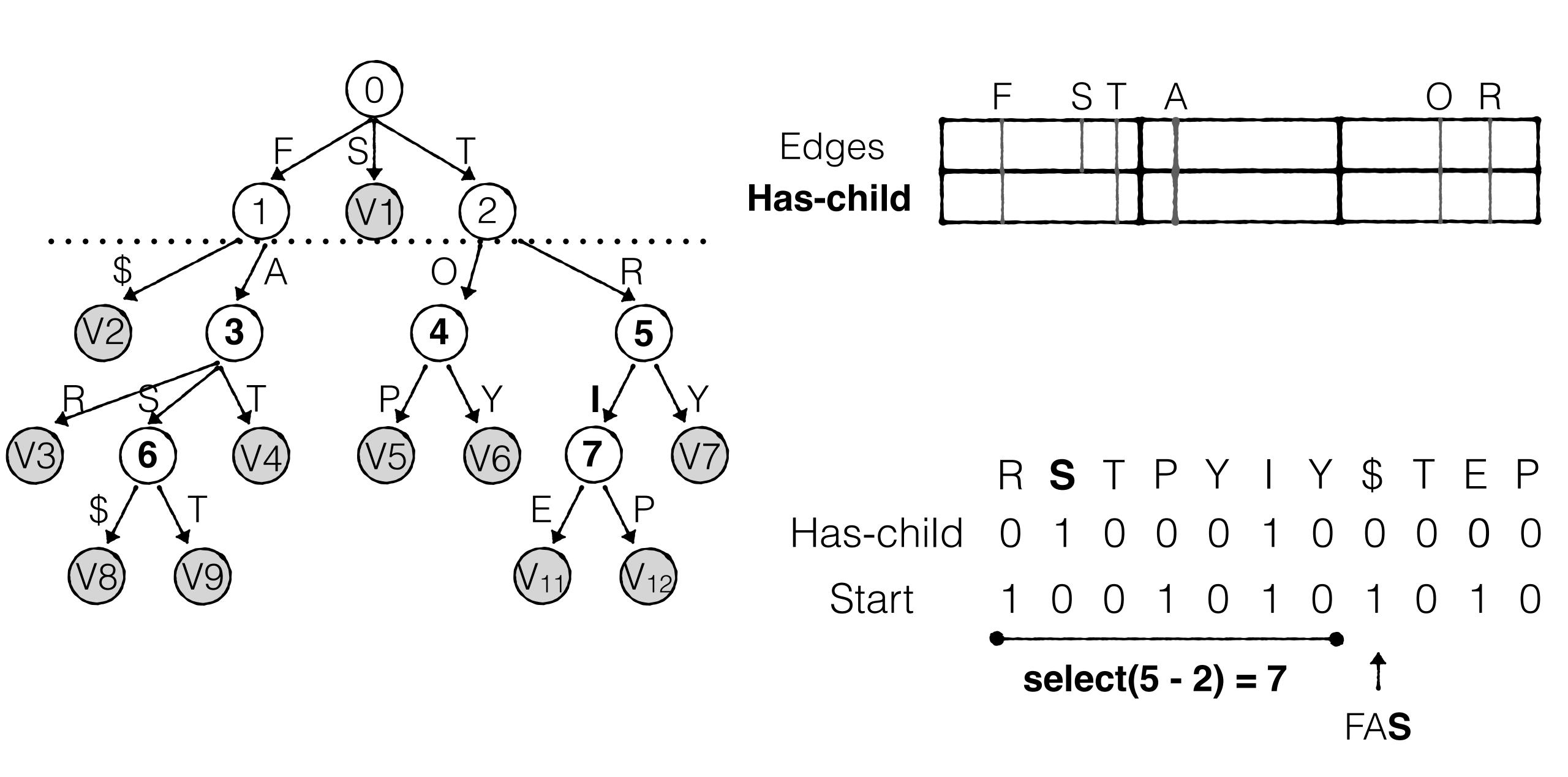


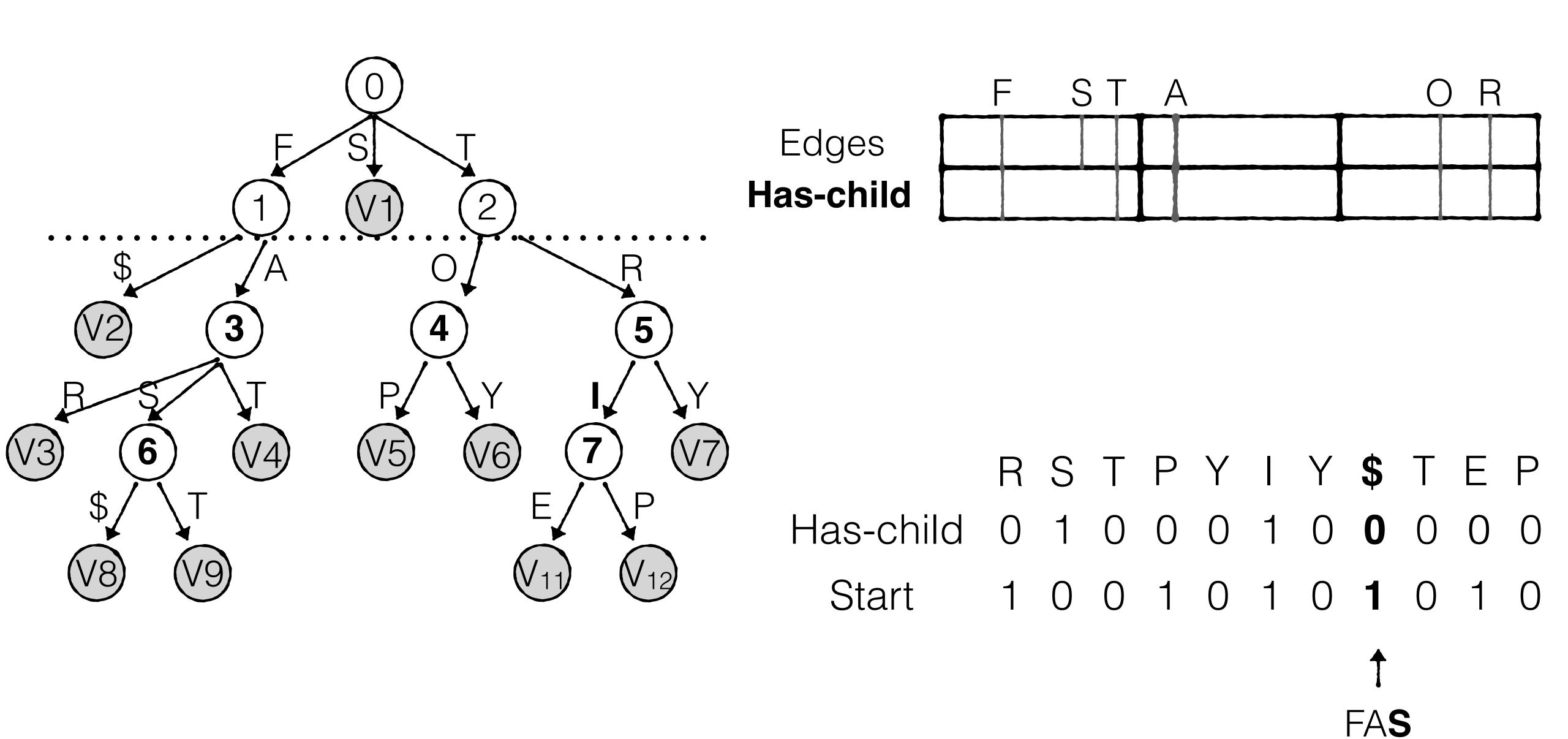


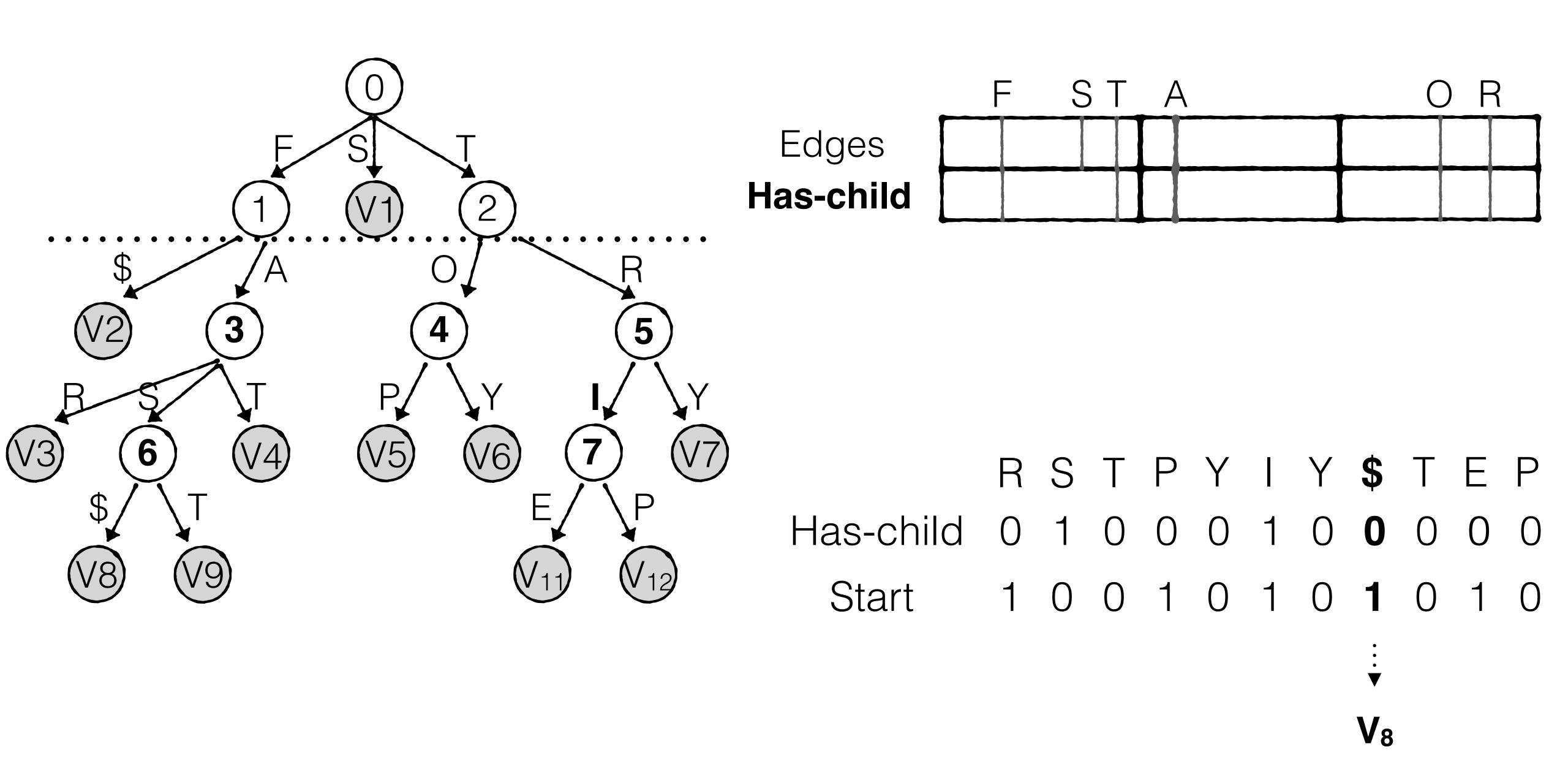


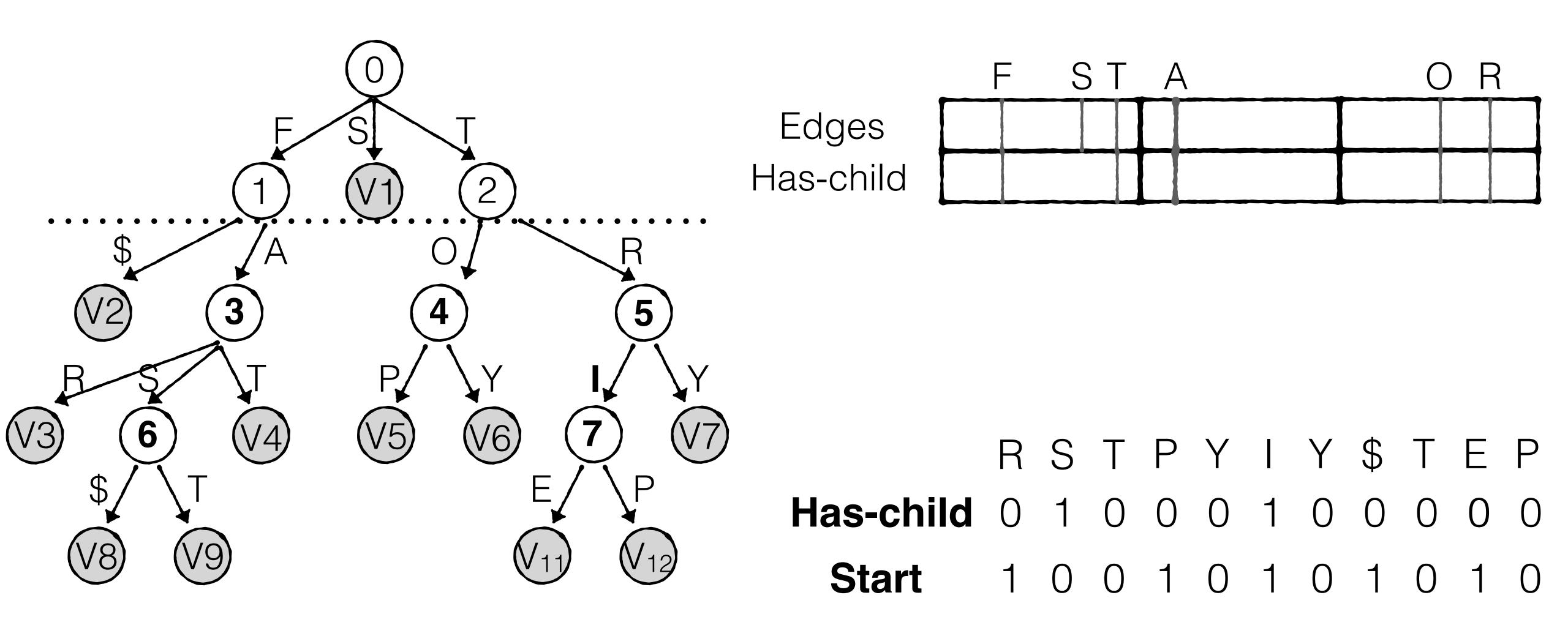




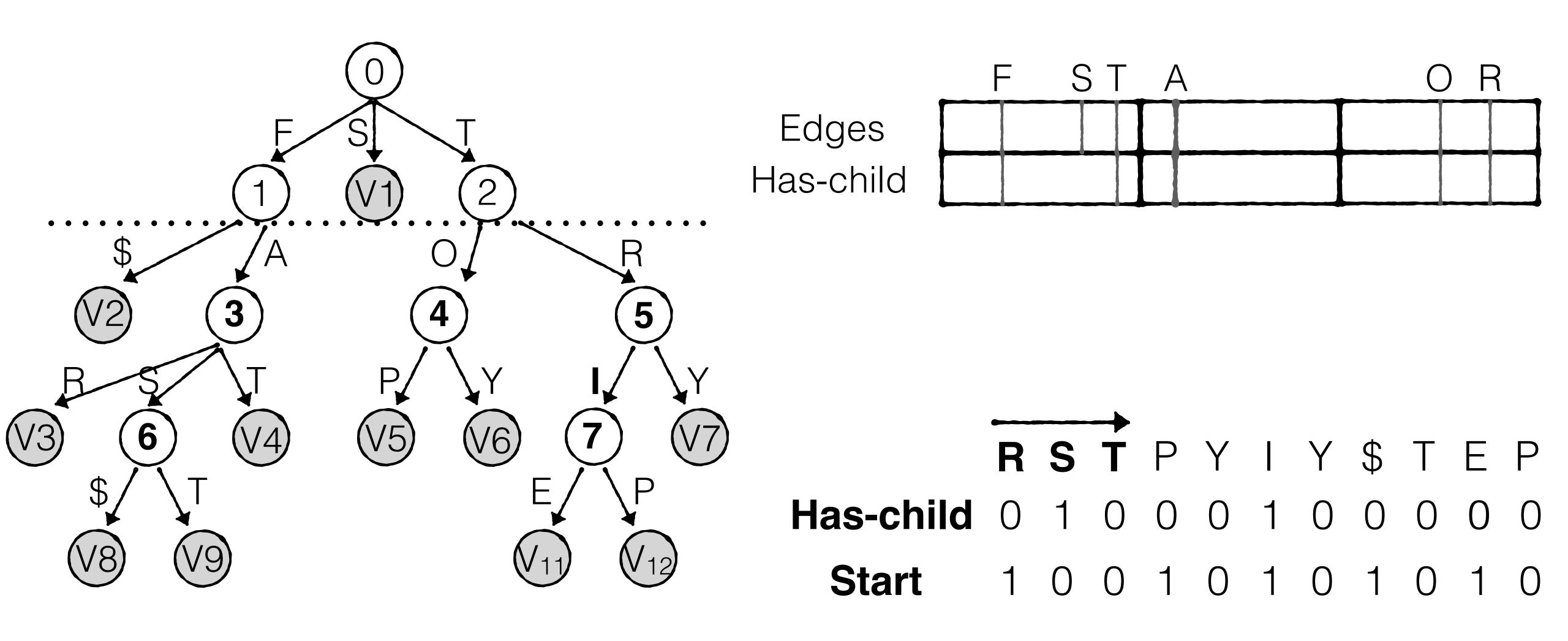








Sparse encoding is space efficient for nodes within many edges



Sparse encoding is space efficient for nodes within many edges

But slower as we must scan edges for each node

## SuRF



Var-length keys & queries

Yes

FPR guarantee

None

Query speed

O(L)

(L = key length)

**Dynamic** 

No

SuRF

Memento

Diva





