

Individual Project Thoughts 1

Functional Dependencies (aka “Determinations” in AI)

Suppose we have a universal relation with attributes A, B, C, ..., each with a set of possible values (e.g., attribute A can have values a1, a2, a3, ... ai)

A	B	C	D	E	F	G	H ...
a1	b3	c2	d5	e7	f3	g1	h6 ...
a4	b2	c2	d4	e2	f1	g1	h5 ...
a2	b1	c1	d2	e5	f5	g3	h2 ...
a1	b3	c3	d5	e6	f4	g1	h8 ...
...							

Suppose we are not told the FDs that are manifest (or intended to be manifest) in this universal relation

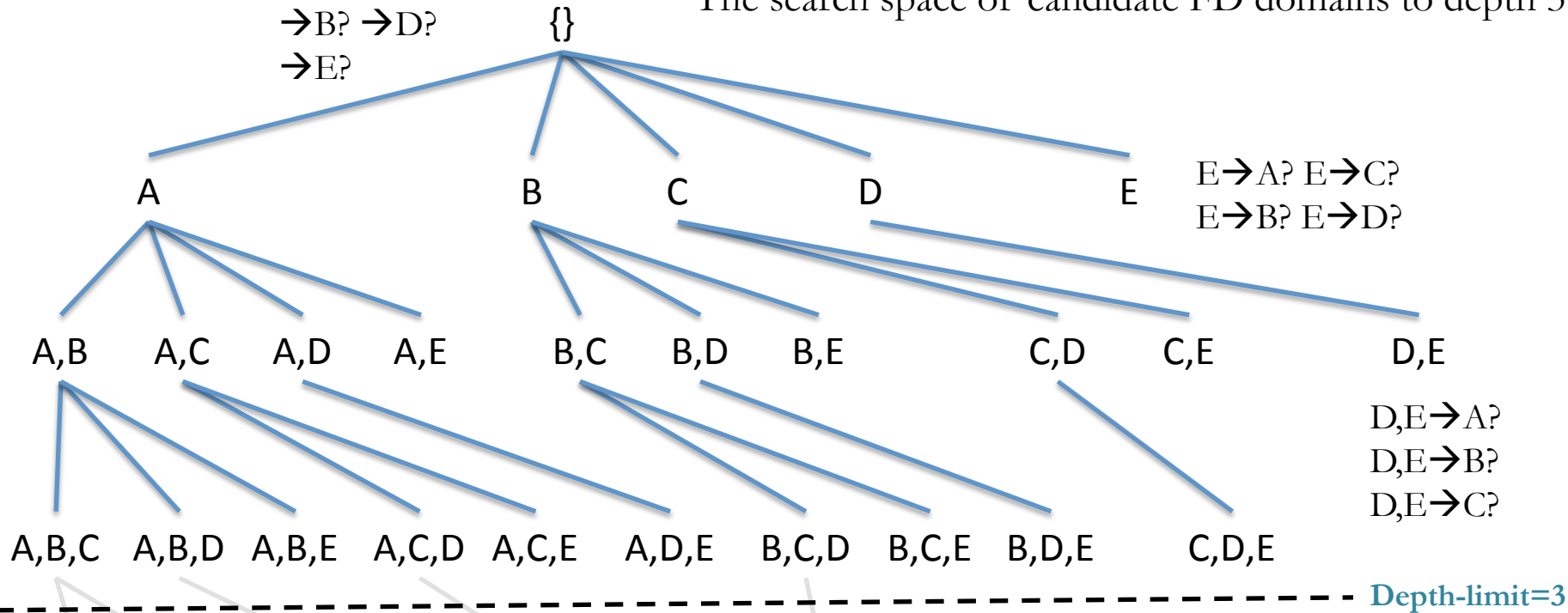
How can we induce the FDs through a process of “unsupervised” machine learning?

Schlimmer, J. (1993). Efficiently Inducing Determinations: A Complete and Systematic Search Algorithm that Uses Optimal Pruning (1993)

<http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.49.2038>

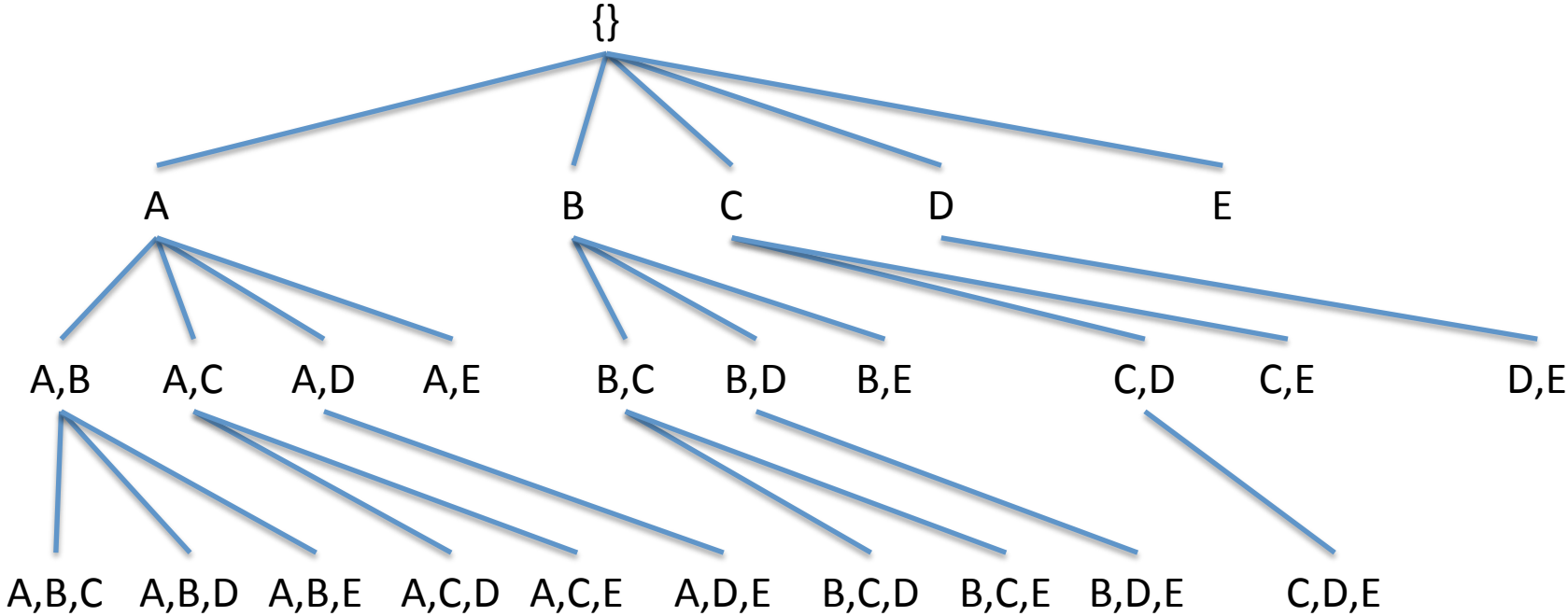
The search space of candidate FD domains to depth 3

$\rightarrow A? \rightarrow C?$
 $\rightarrow B? \rightarrow D?$
 $\rightarrow E?$



$A, B, C \rightarrow D?$	$A, B, E \rightarrow C?$	$A, C, E \rightarrow B?$	$B, C, D \rightarrow A?$	$B, D, E \rightarrow A?$	$C, D, E \rightarrow A?$
$A, B, C \rightarrow E?$	$A, B, E \rightarrow D?$	$A, C, E \rightarrow D?$	$B, C, D \rightarrow E?$	$B, D, E \rightarrow C?$	$C, D, E \rightarrow B?$
$A, B, D \rightarrow C?$	$A, C, D \rightarrow B?$	$A, D, E \rightarrow B?$	$B, C, E \rightarrow A?$		
$A, B, D \rightarrow E?$	$A, C, D \rightarrow E?$	$A, D, E \rightarrow C?$	$B, C, E \rightarrow D?$		

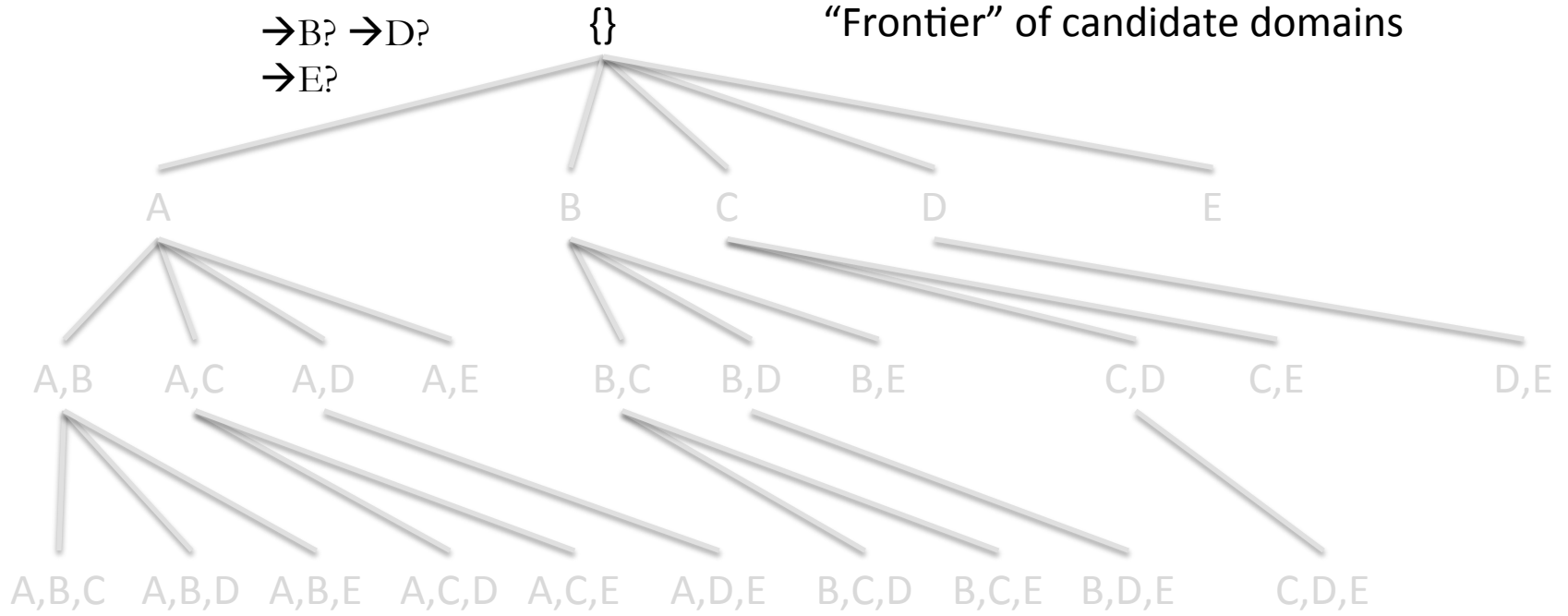
This space can be searched by breadth-first search, depth-first search, or a more sophisticated heuristic search



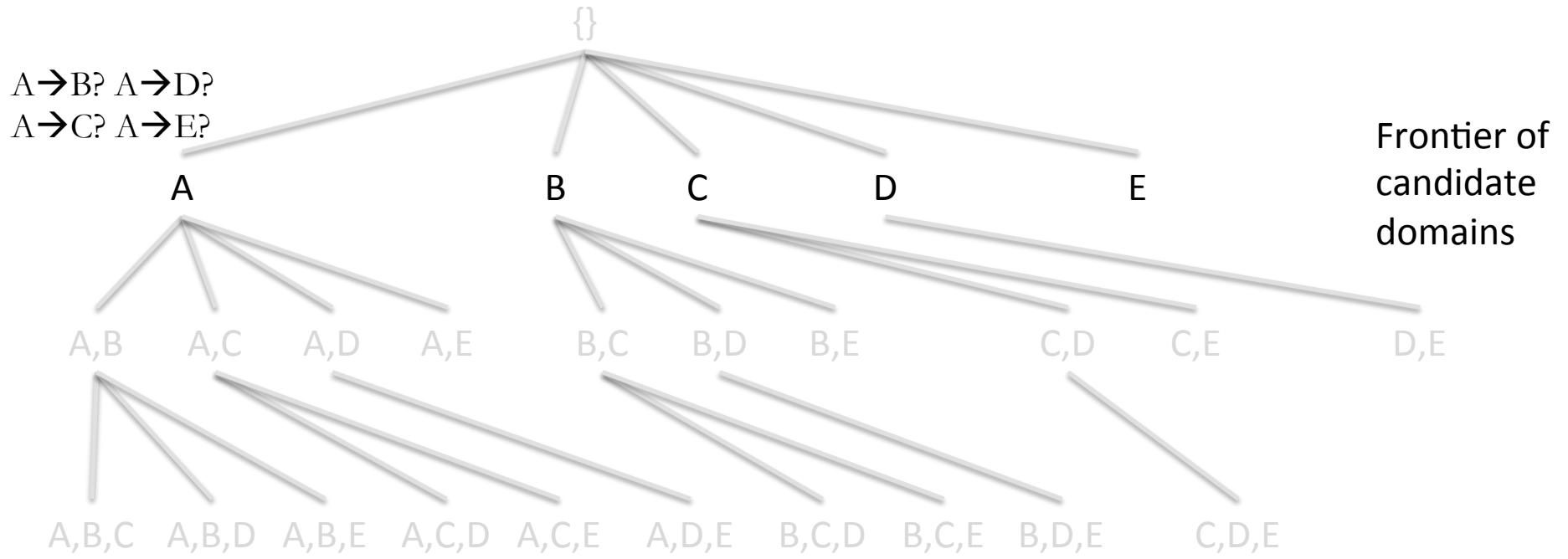
Breadth first search

→A? →C?
→B? →D?
→E?

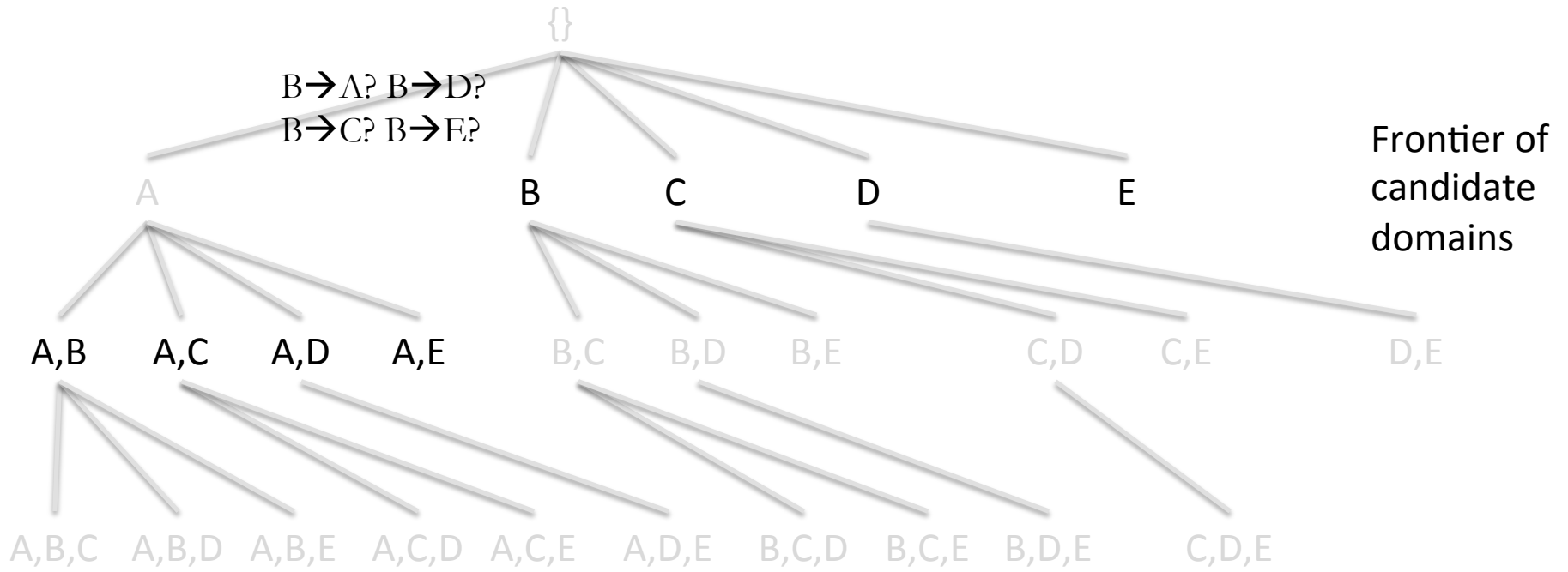
“Frontier” of candidate domains



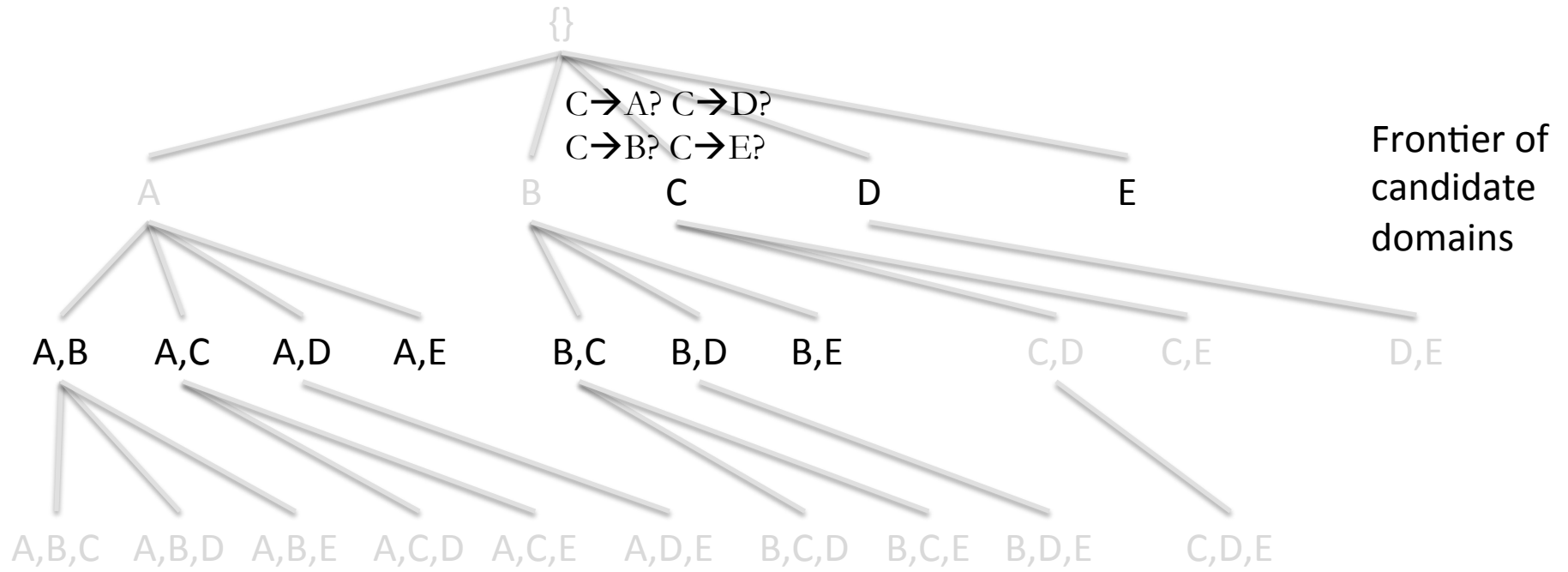
Breadth first search



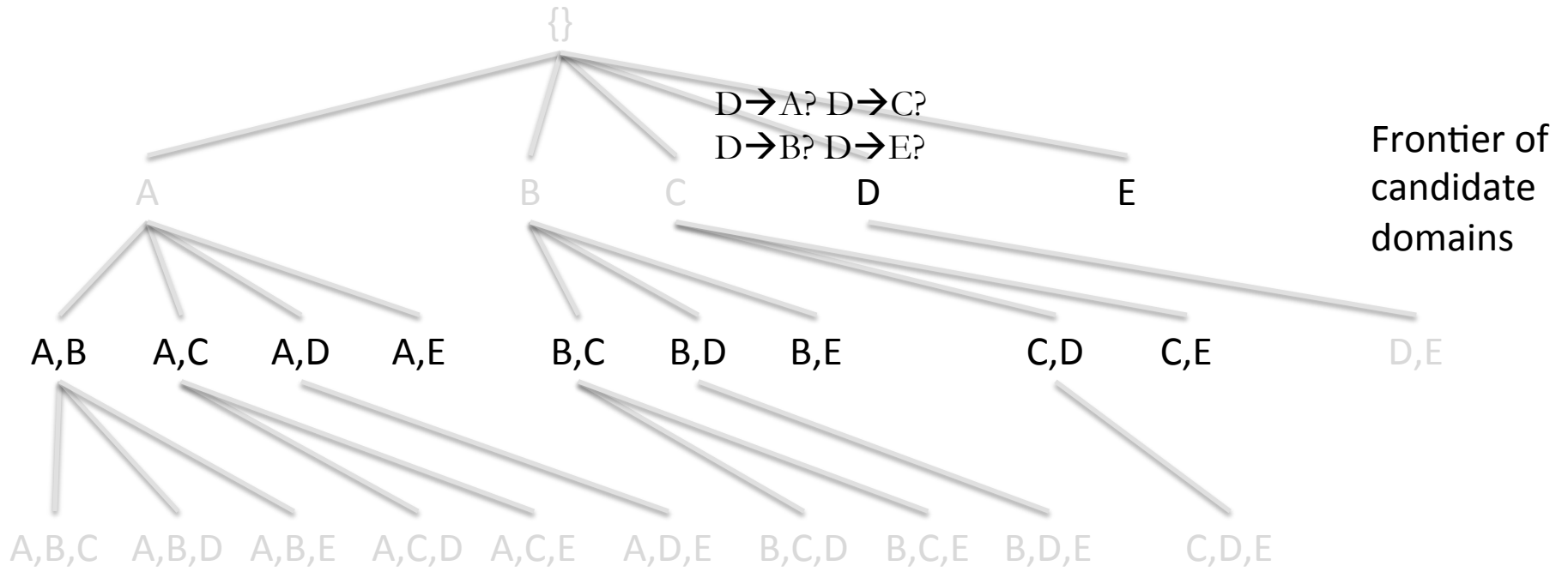
Breadth first search



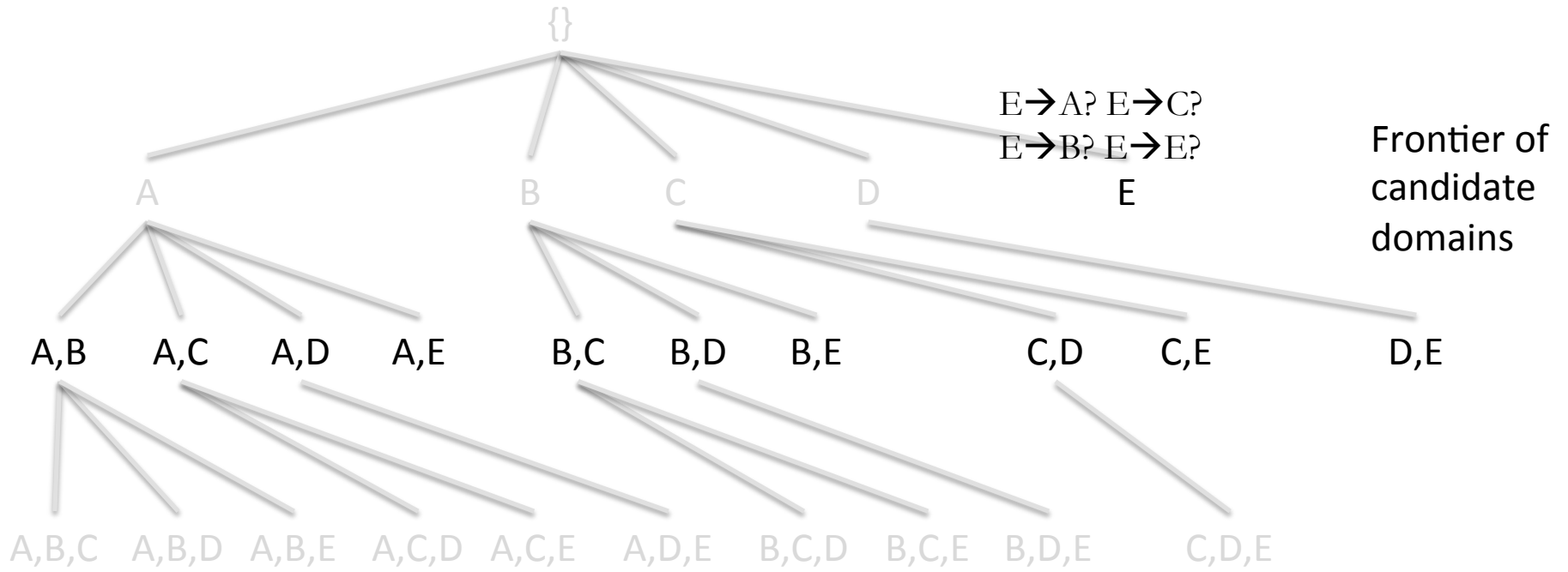
Breadth first search



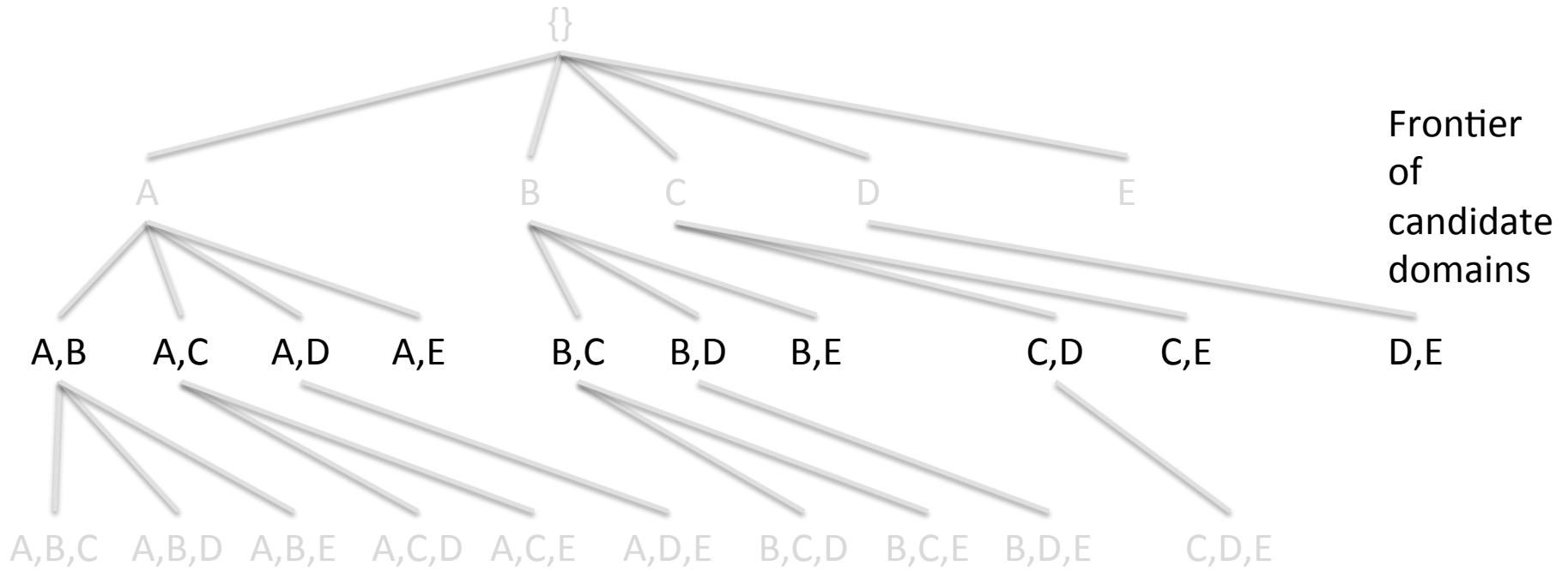
Breadth first search



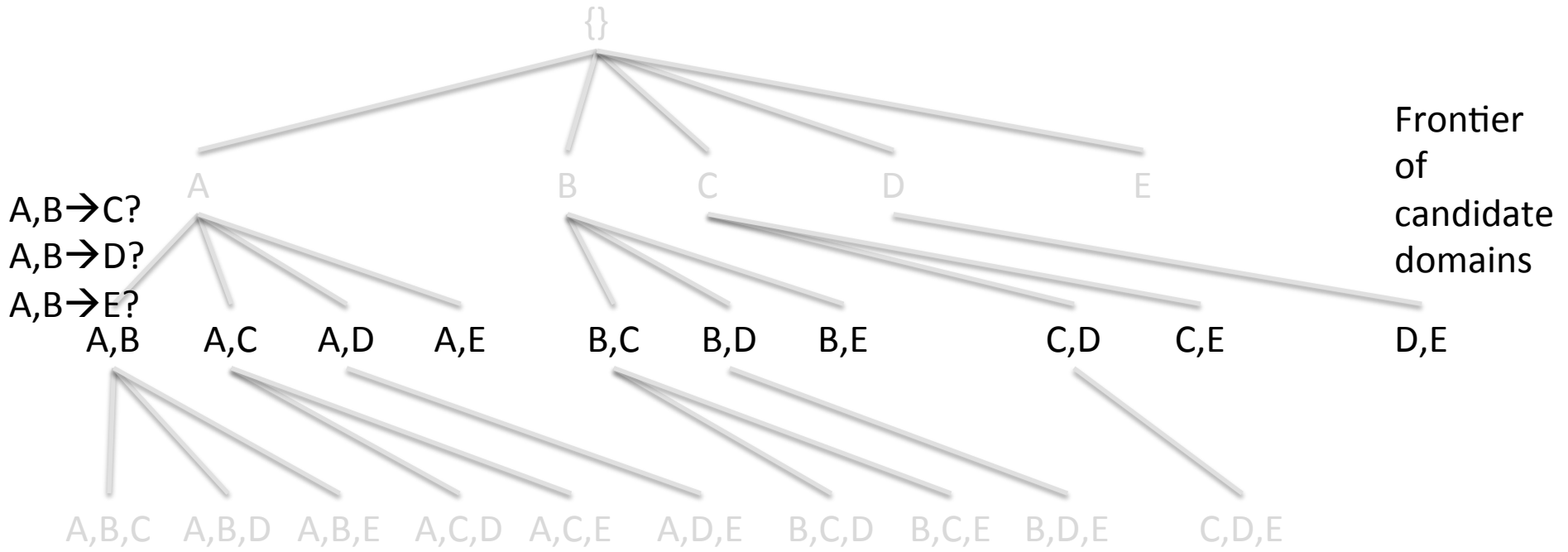
Breadth first search



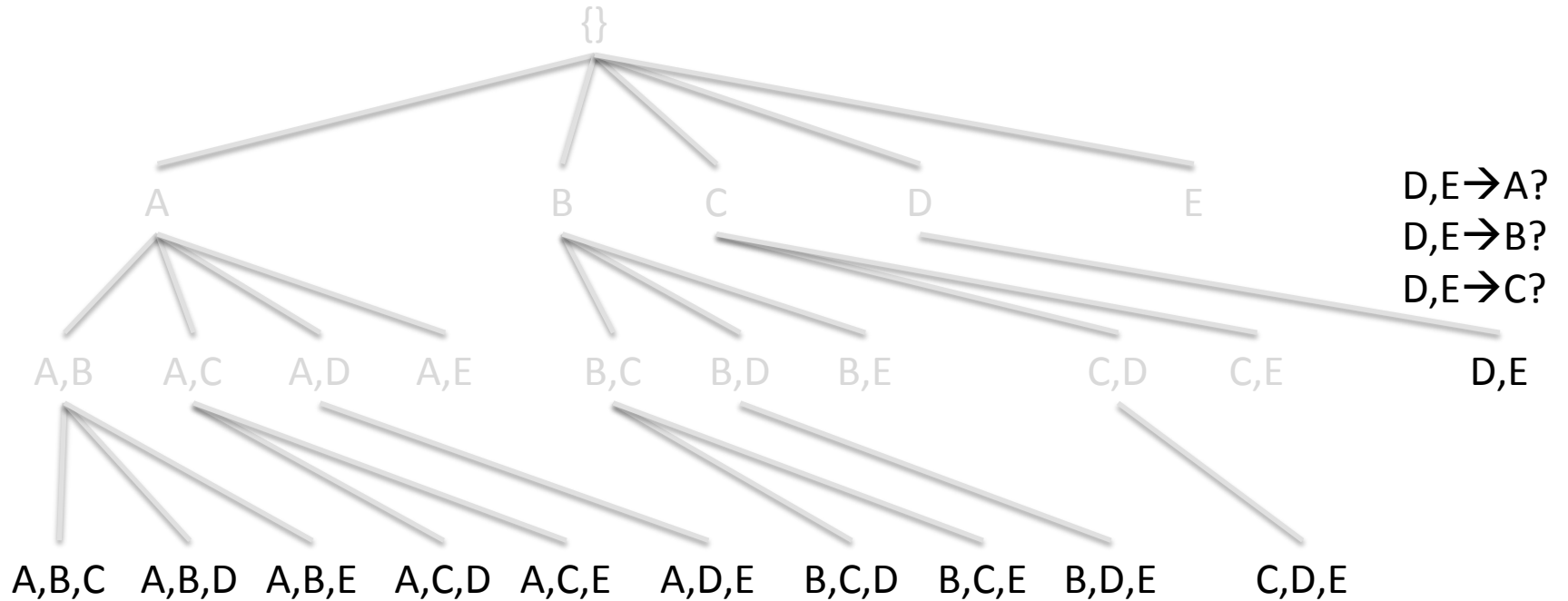
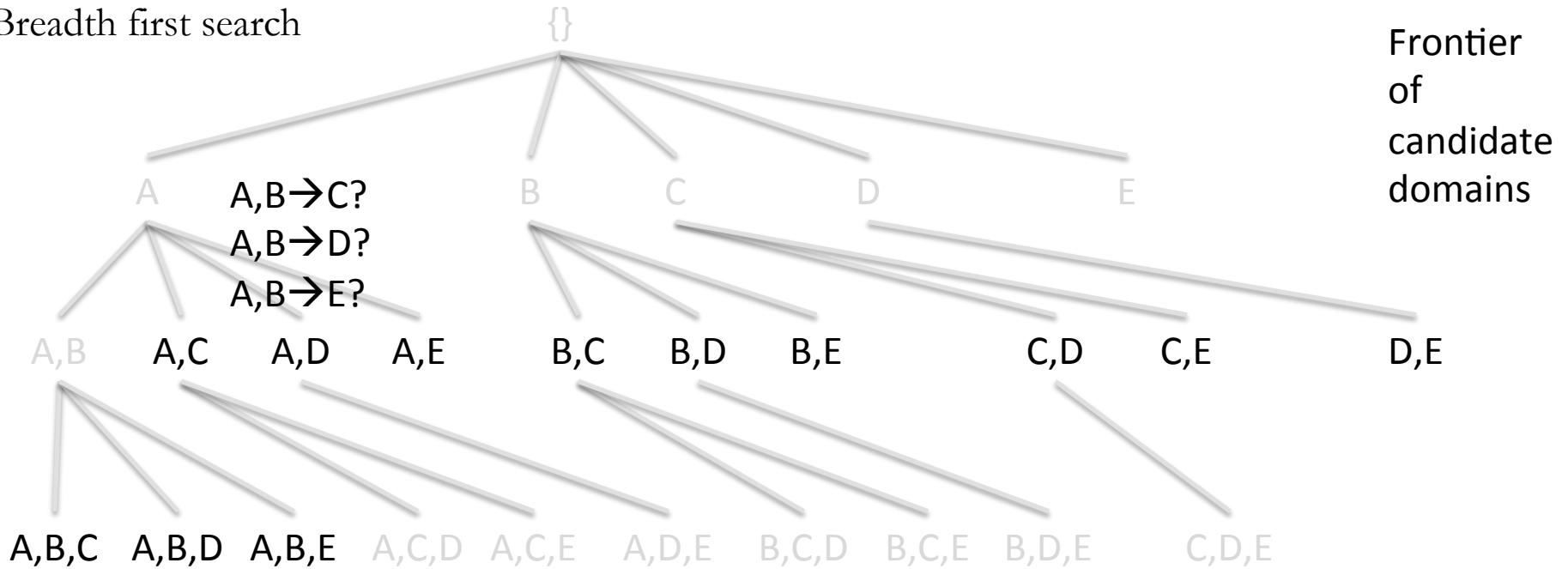
Breadth first search



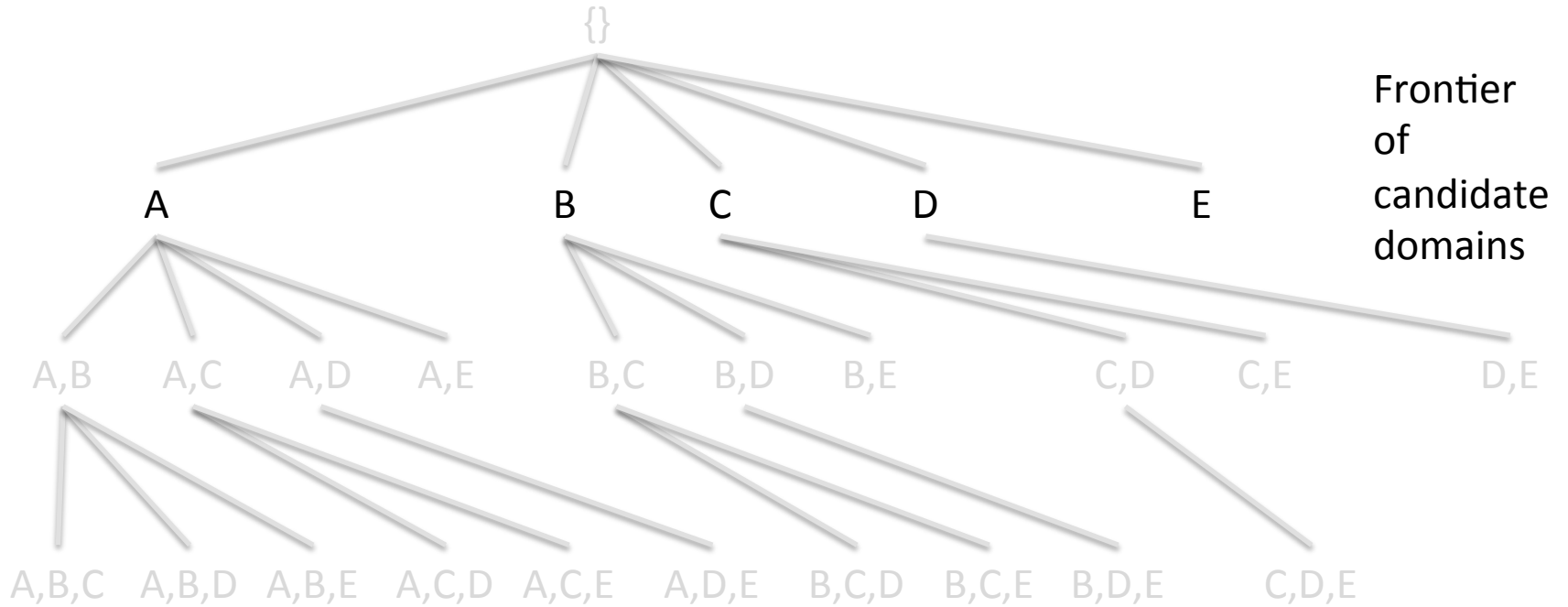
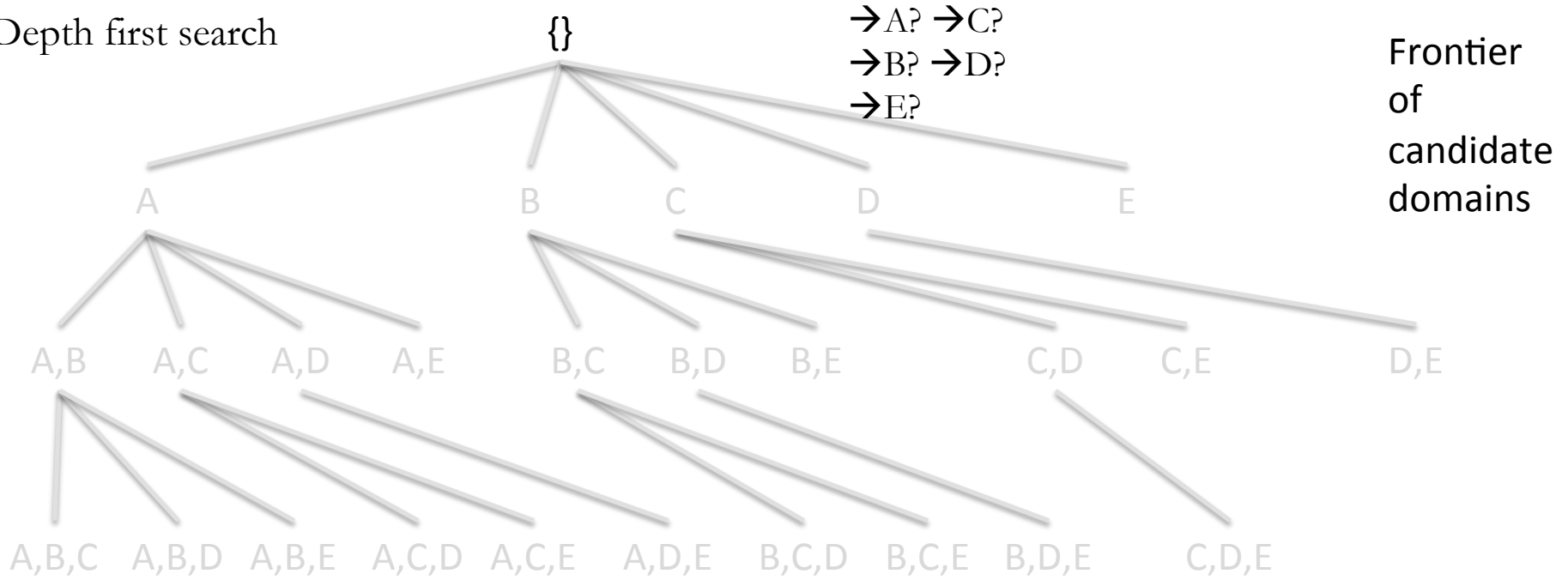
Breadth first search



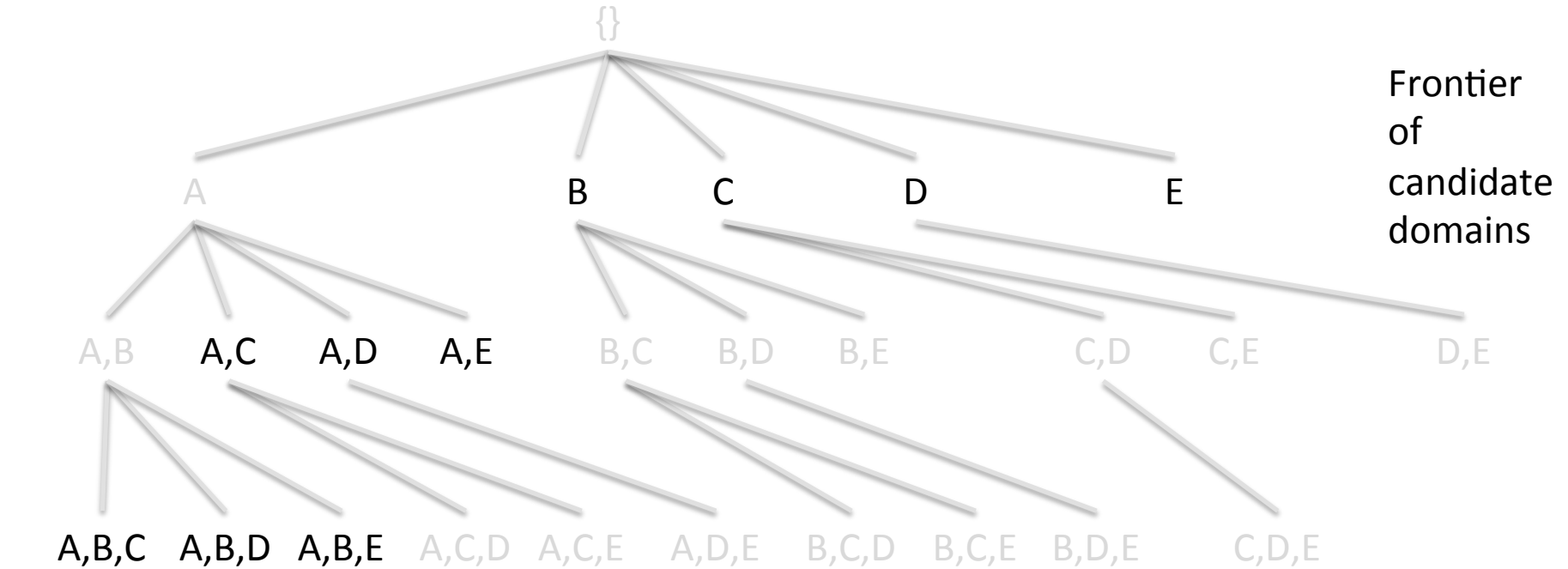
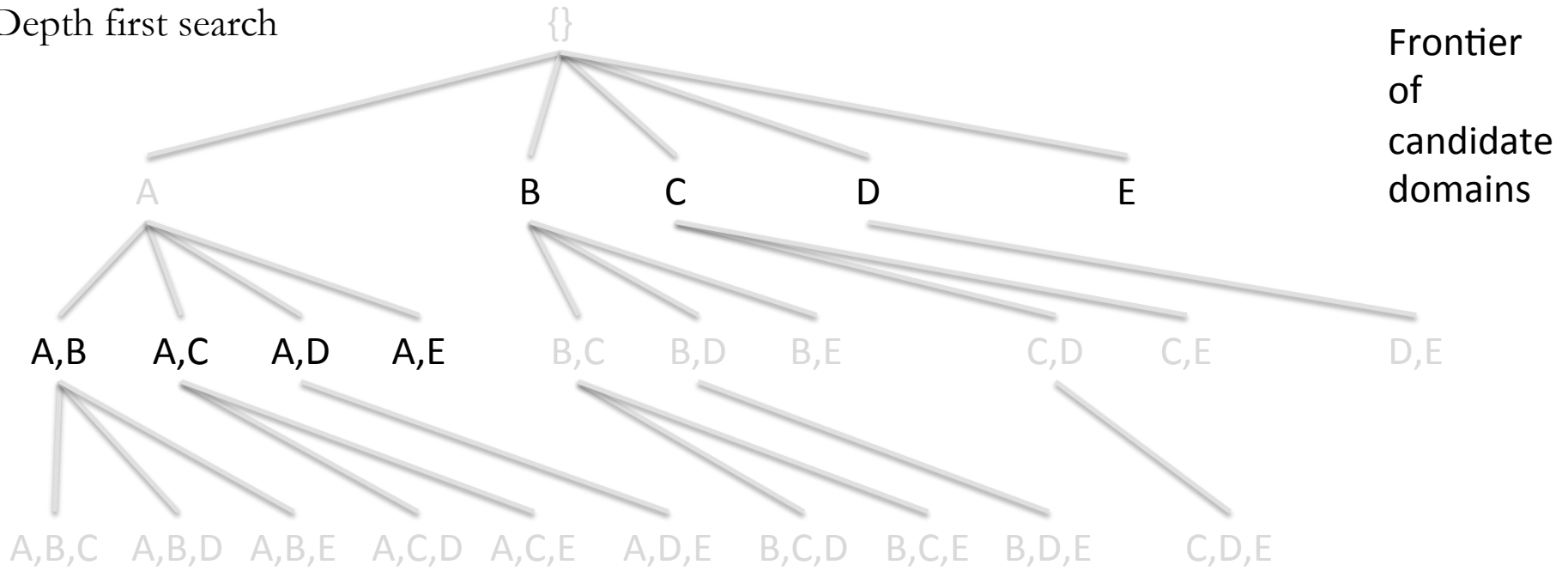
Breadth first search



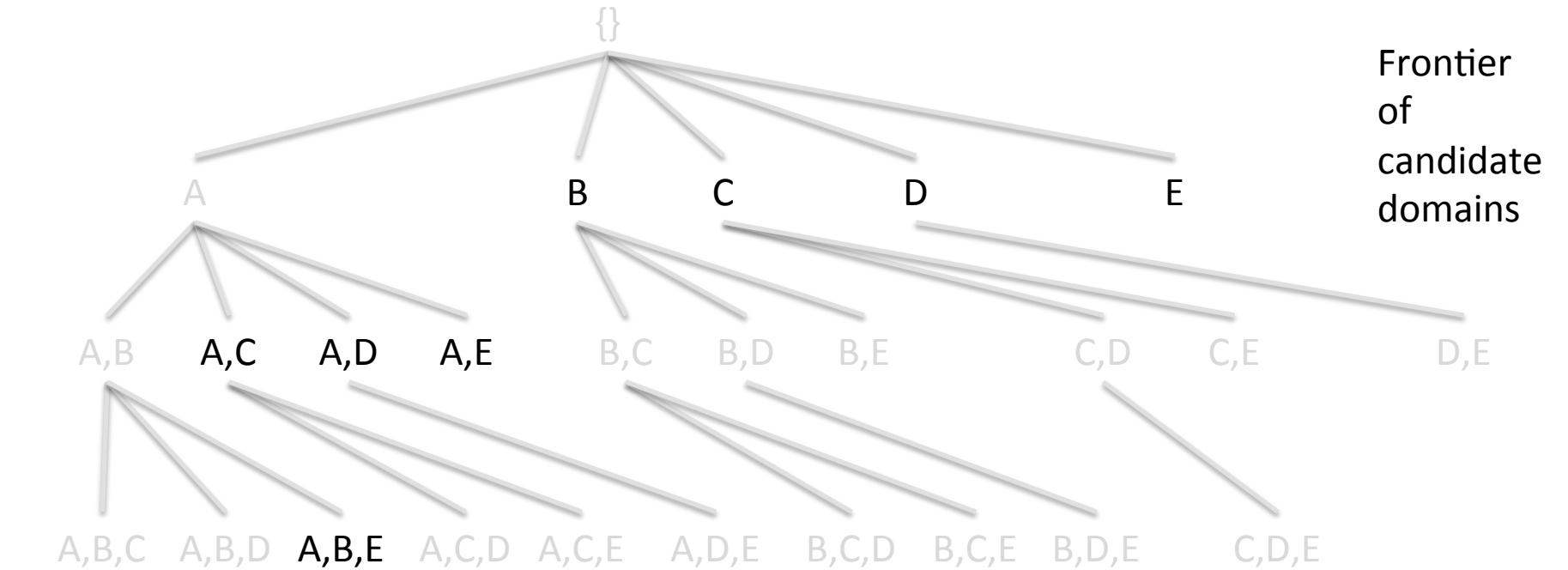
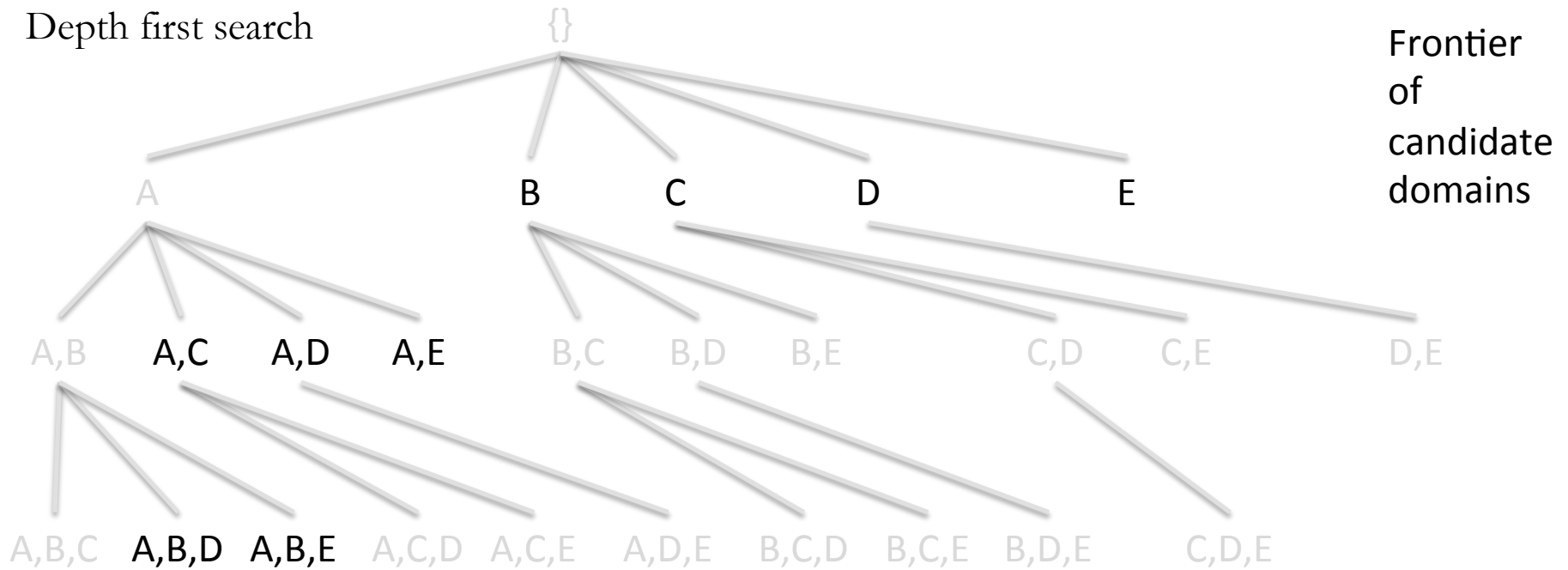
Depth first search



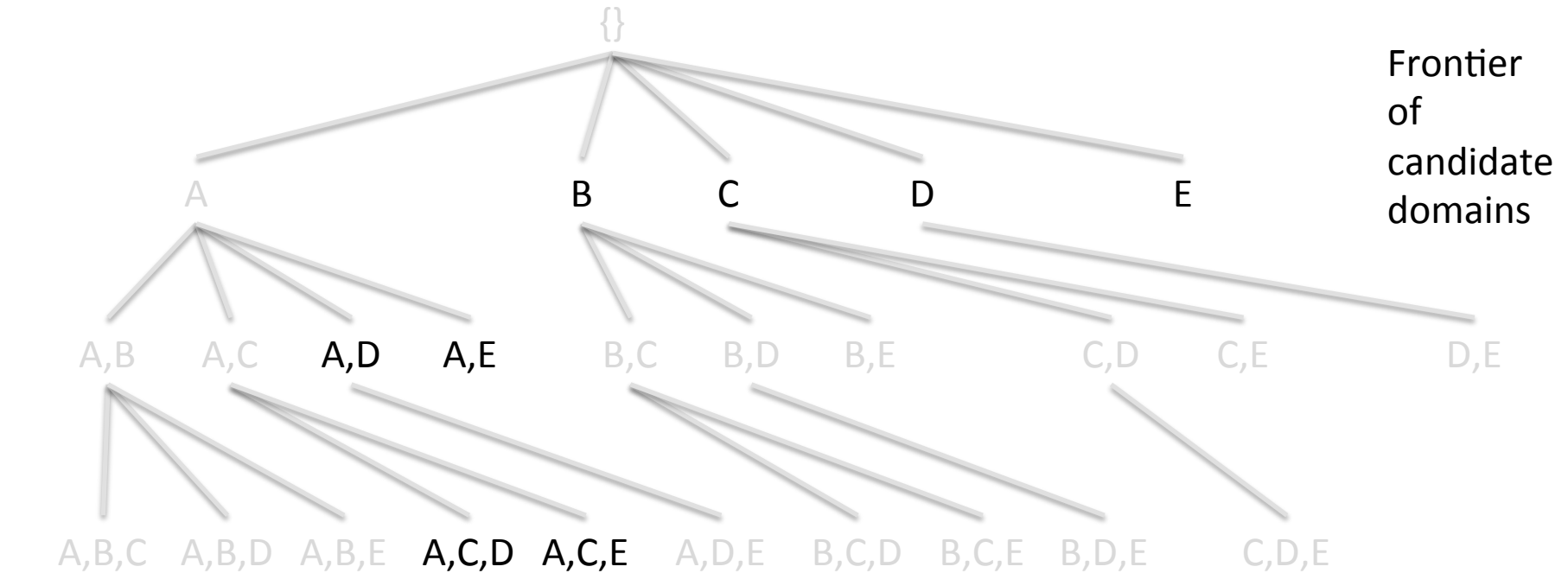
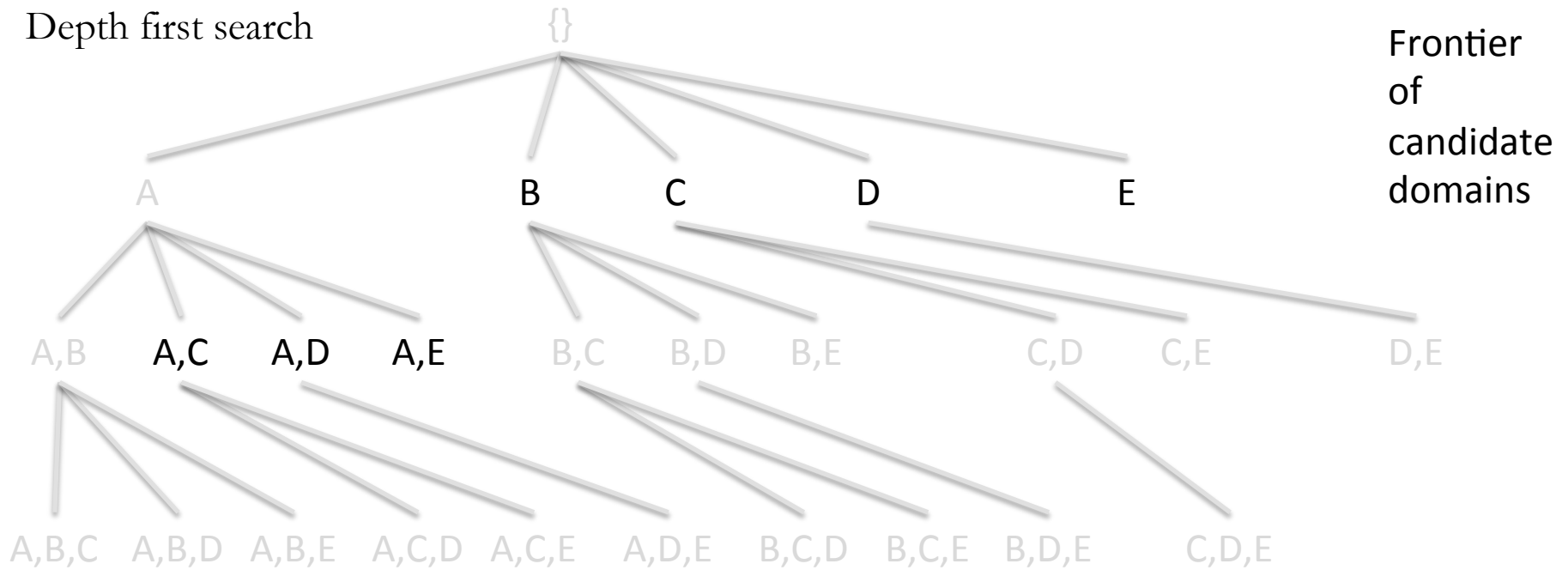
Depth first search



Depth first search

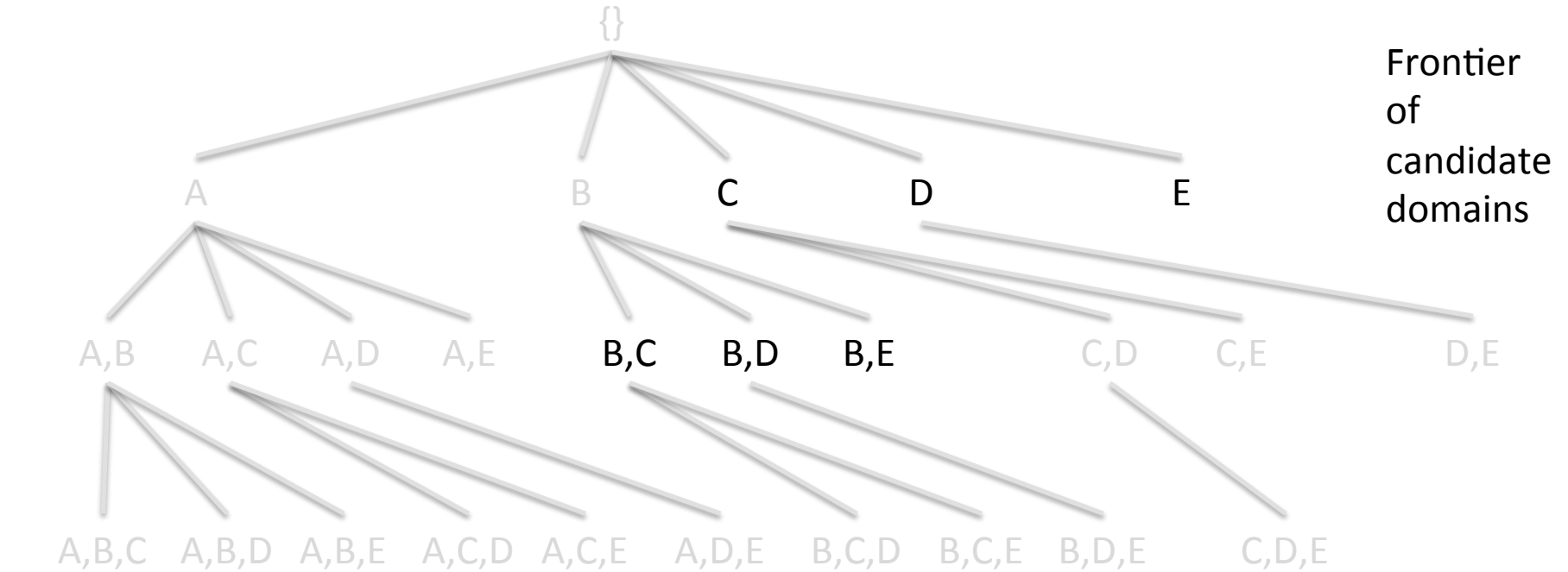
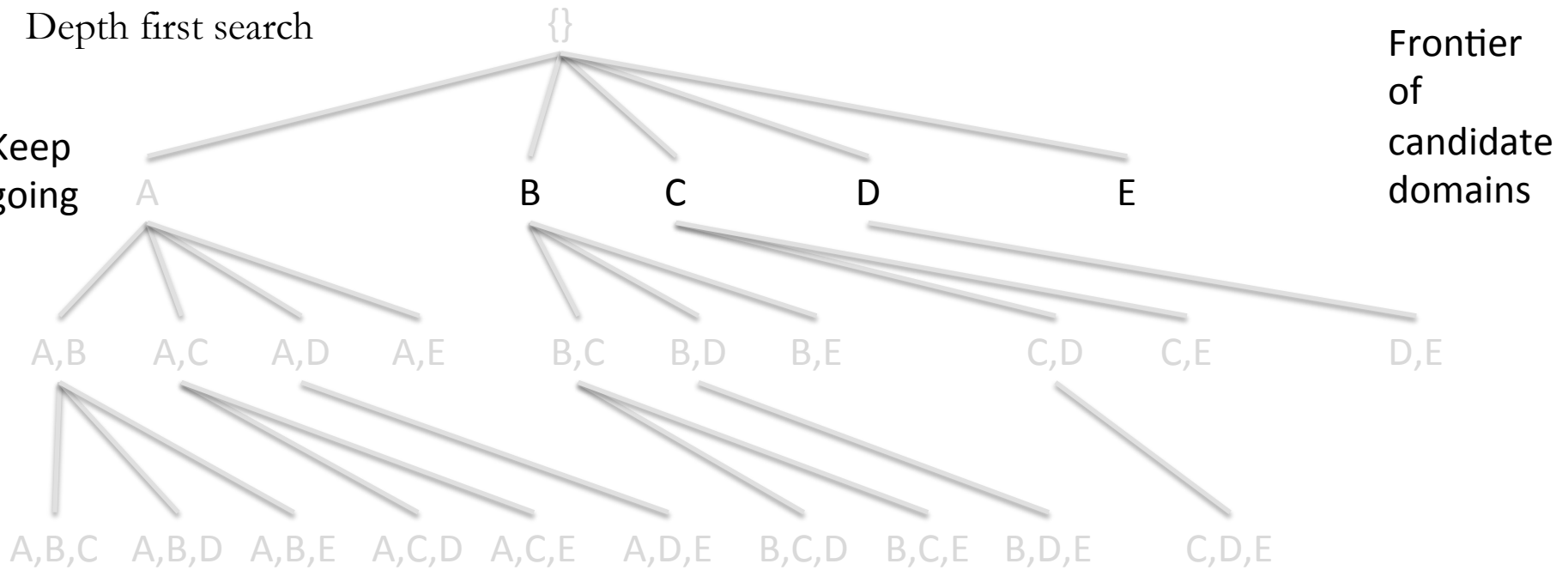


Depth first search



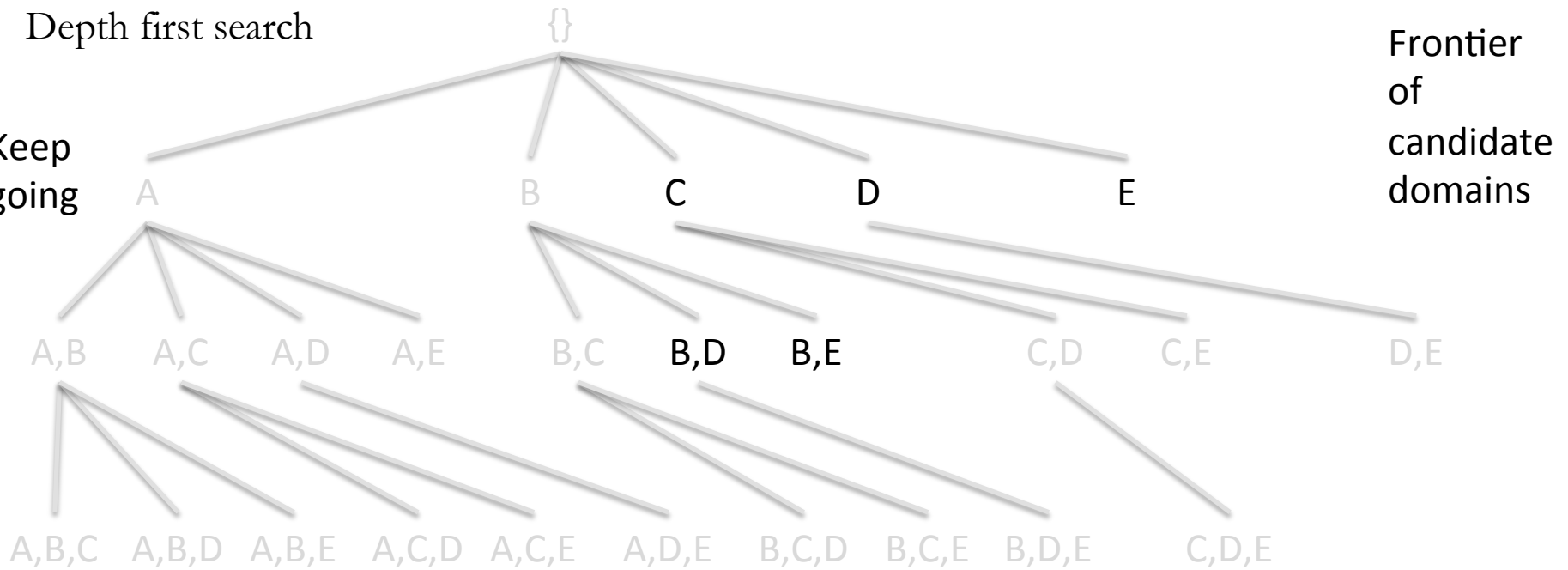
Depth first search

Keep going

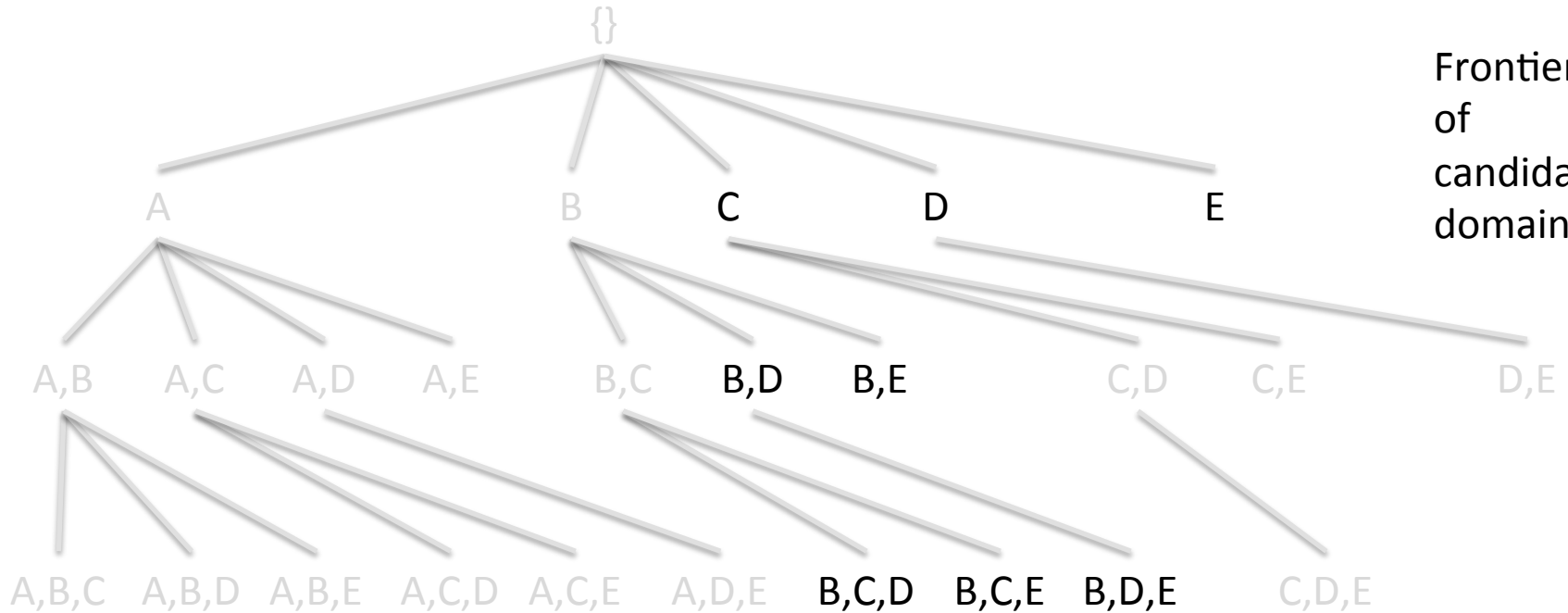


Depth first search

Keep going

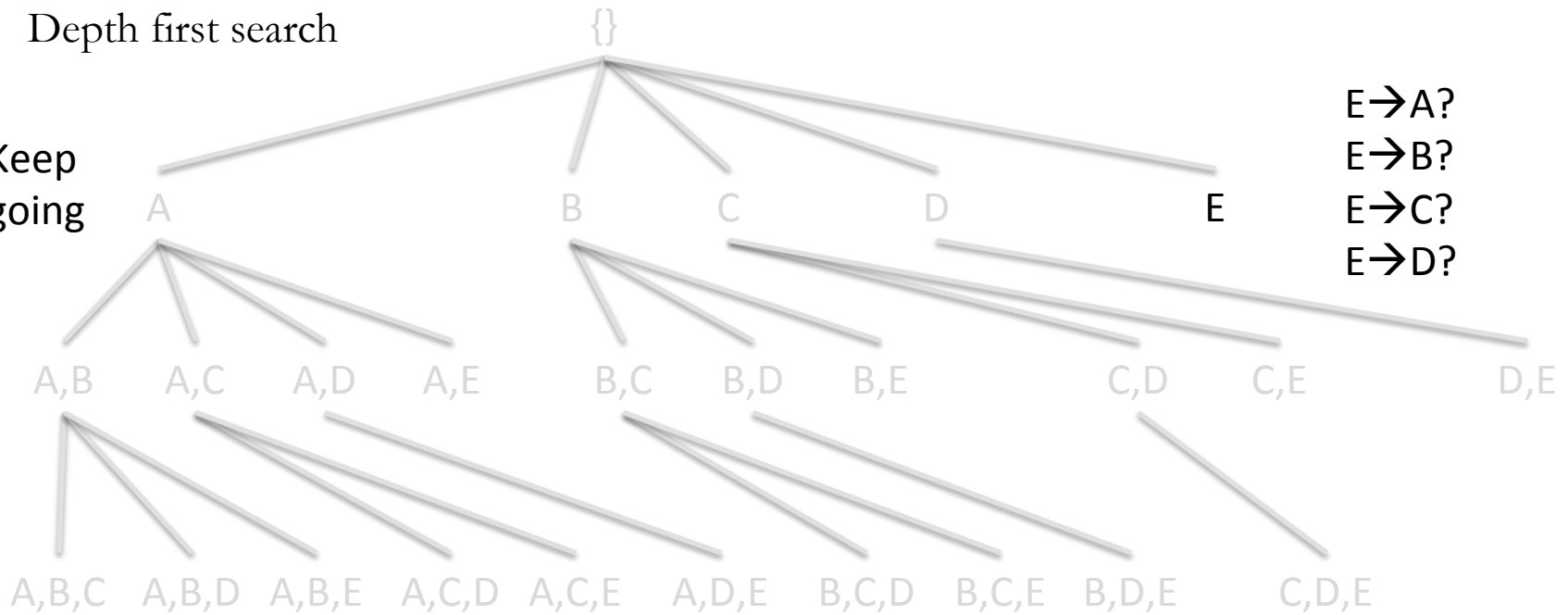


Frontier of candidate domains

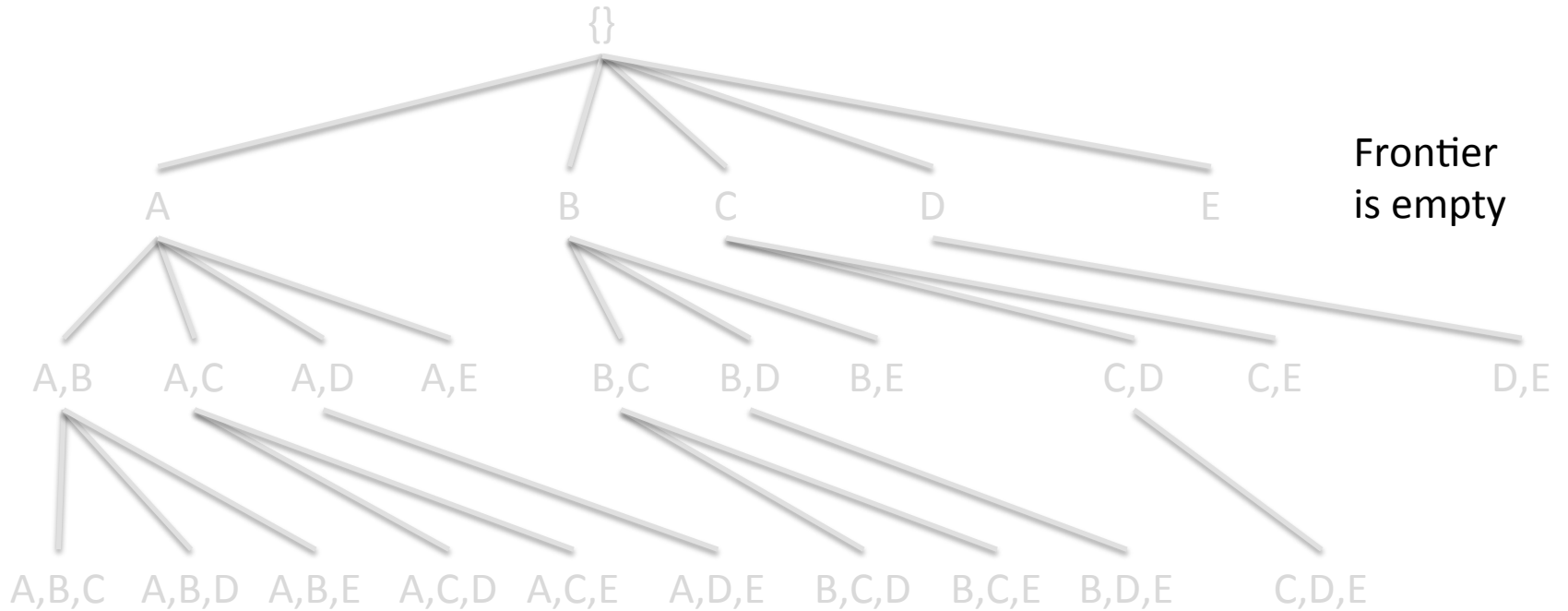


Depth first search

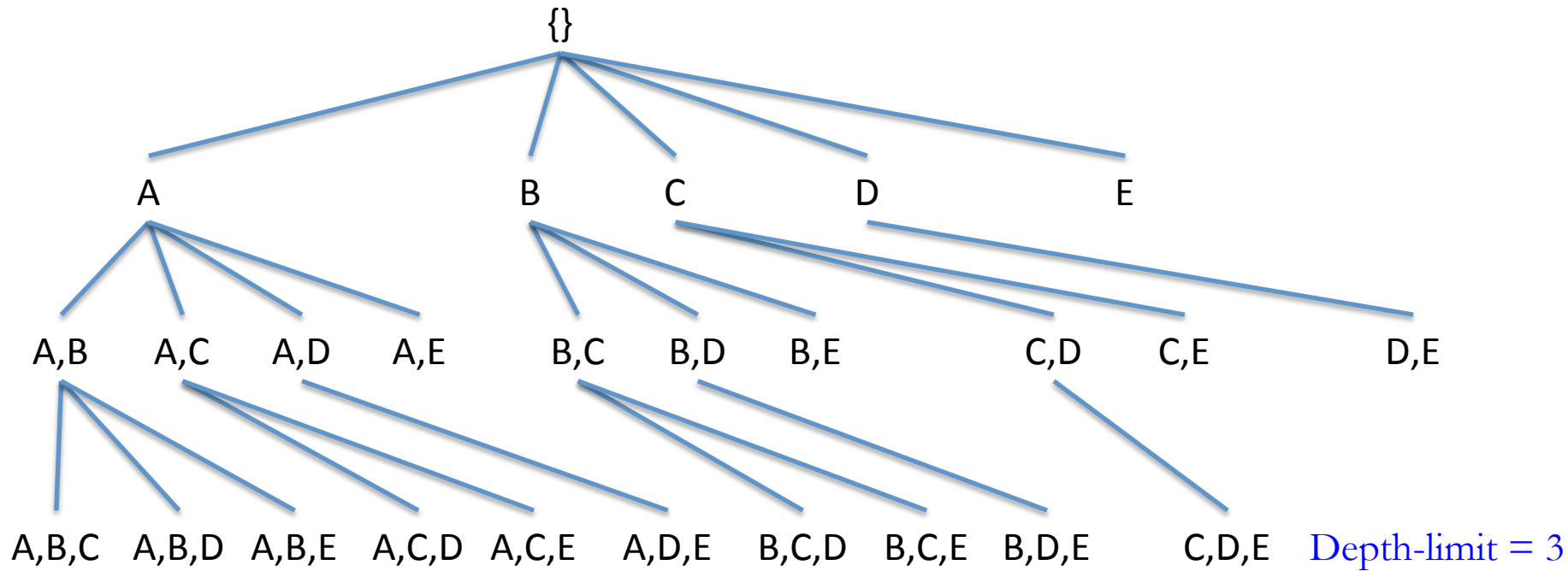
Keep going



Frontier is empty



Instead of learning only perfectly consistent FDs, beneficial to learn approximate FDs (almost perfectly consistent) This



A,B,C \rightarrow D?

((a1,b3,c4:100 rows),((d1:98 rows), (d3: 2 rows))
 ((a2,b2,c1:43 rows),((d2:42 rows)(d1: 1 row))
 ((a3,b1,c1:15 rows),((d1:15 rows))

Number of rows adhering to dominant pattern with Range attribute (numerator)

Number of rows adhering to Domain patterns(denominator)

$$\frac{(98+42+15)}{(100+43+15)} = 155/158 = 0.98 \text{ support}$$

If parameter *minimal-support* = 0.95 then accept A,B,C \rightarrow D (0.98)

Pseudo code

```
Find-Approximate-Functional-Dependencies (data-set, depth-limit, minimal-support)
  approximate-FDs  $\leftarrow$   $\{\}$  /* the empty set */
  domains-frontier  $\leftarrow$   $\{\{\}\}$  /* the set containing the empty set */
  WHILE domains-frontier  $\neq$   $\{\}$ 
    next-domain  $\leftarrow$  Select from domains-frontier
    domains-frontier  $\leftarrow$  domains-frontier  $-$  next-domain
    FOR each attribute, Y, where Y  $\neq$  any attribute in next-domain
      support  $\leftarrow$  compute support for FD (next-domain  $\rightarrow$  Y) using data-set
      IF support  $\geq$  minimum-support
        THEN approximate-FDs  $\leftarrow$  approximate-FDs + (next-domain, Y, support)
      IF  $|$ next-domain $|$   $<$  depth-limit
        THEN FOR all attributes, X, where X  $>$  all attributes in next-domain
          domains-frontier  $\leftarrow$  domains-frontier + (next-domain + X)
  RETURN approximate-FDs
```

Pseudo code

Find-Approximate-Functional-Dependencies (*data-set*, *depth-limit*, *minimal-support*)

approximate-FDs \leftarrow $\{\}$ /* the empty set */

domains-frontier \leftarrow $\{\{\}\}$ /* the set containing the empty set */

WHILE *domains-frontier* \neq $\{\}$

next-domain \leftarrow Select from *domains-frontier*

domains-frontier \leftarrow *domains-frontier* $-$ *next-domain*

FOR each attribute, *Y*, where *Y* \neq any attribute in *next-domain*

support \leftarrow compute support for FD (*next-domain* \rightarrow *Y*) using *data-set*

 IF *support* \geq *minimum-support*

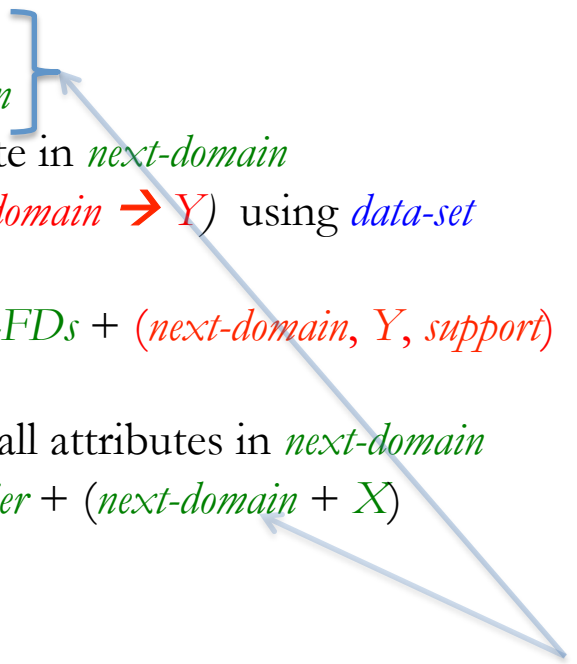
 THEN *approximate-FDs* \leftarrow *approximate-FDs* + (*next-domain*, *Y*, *support*)

IF $|$ *next-domain* $|$ $<$ *depth-limit*

 THEN FOR all attributes, *X*, where *X* $>$ all attributes in *next-domain*

domains-frontier \leftarrow *domains-frontier* + (*next-domain* + *X*)

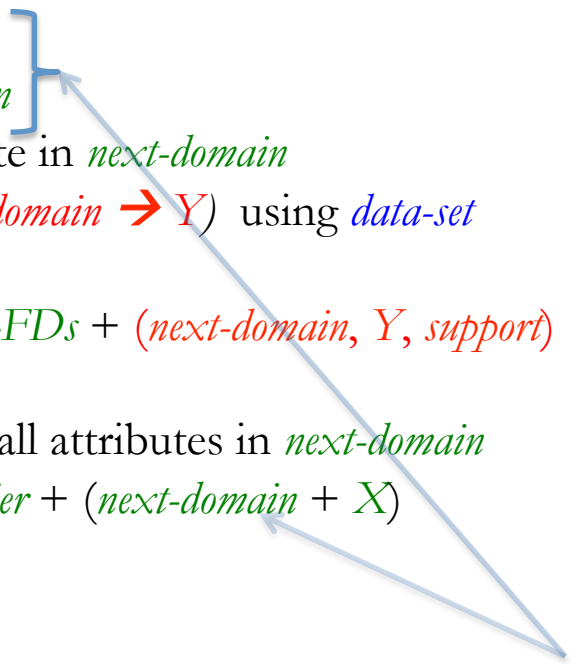
RETURN *approximate-FDs*



If *domains-frontier*
is a queue
then breadth-first
search

Pseudo code

```
Find-Approximate-Functional-Dependencies (data-set, depth-limit, minimal-support)
  approximate-FDs  $\leftarrow$  { } /* the empty set */
  domains-frontier  $\leftarrow$  { { } } /* the set containing the empty set */
  WHILE domains-frontier  $\neq$  { }
    next-domain  $\leftarrow$  Select from domains-frontier
    domains-frontier  $\leftarrow$  domains-frontier - next-domain
  FOR each attribute, Y, where Y  $\neq$  any attribute in next-domain
    support  $\leftarrow$  compute support for FD (next-domain  $\rightarrow$  Y) using data-set
    IF support  $\geq$  minimum-support
      THEN approximate-FDs  $\leftarrow$  approximate-FDs + (next-domain, Y, support)
    IF |next-domain| < depth-limit
      THEN FOR all attributes, X, where X  $\supset$  all attributes in next-domain
        domains-frontier  $\leftarrow$  domains-frontier + (next-domain + X)
  RETURN approximate-FDs
```



If domains-frontier
is a stack
then depth-first
search