

Lecture 10

Finding strongly connected components

Announcement: Treehacks!

- [Special guest announcers...]

Announcements

- HW4 due Friday
- Nothing assigned Friday because...
- **MIDTERM** in class, Wednesday 2/20.
 - During class, 10:30-11:50
 - Locations TBD
 - You may bring one double-sided letter-size page of notes, that *you have prepared yourself.*
 - Otherwise it is closed-book.
 - Any material through BFS/DFS (Lecture 9) is fair game.
 - Practice exams on the website.

Last time

- Breadth-first and depth-first search
- Plus, applications!
 - Topological sorting
 - Clarification: Does it work if you don't start at a source?
 - Answer: It does!! Try it ☺
 - In-order traversal of BSTs
 - Shortest path in unweighted graphs
 - Testing bipartite-ness
- The key was paying attention to the structure of the tree that these search algorithms implicitly build.

Today

- One more application:

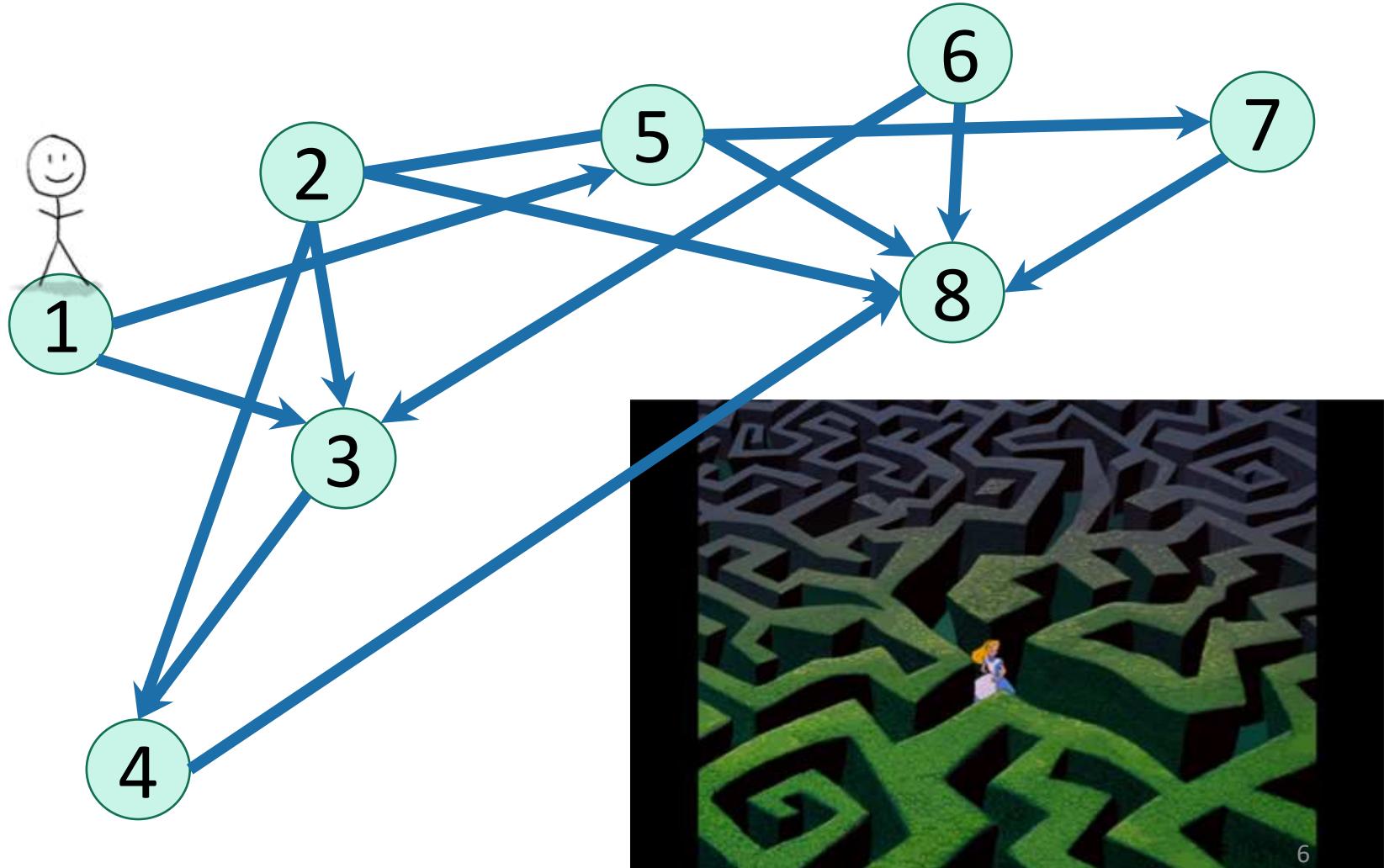
Finding **Strongly Connected Components**

- But first! Let's briefly recap DFS...

Today, all graphs are **directed**!
Check that the things we did
on Monday still all work!

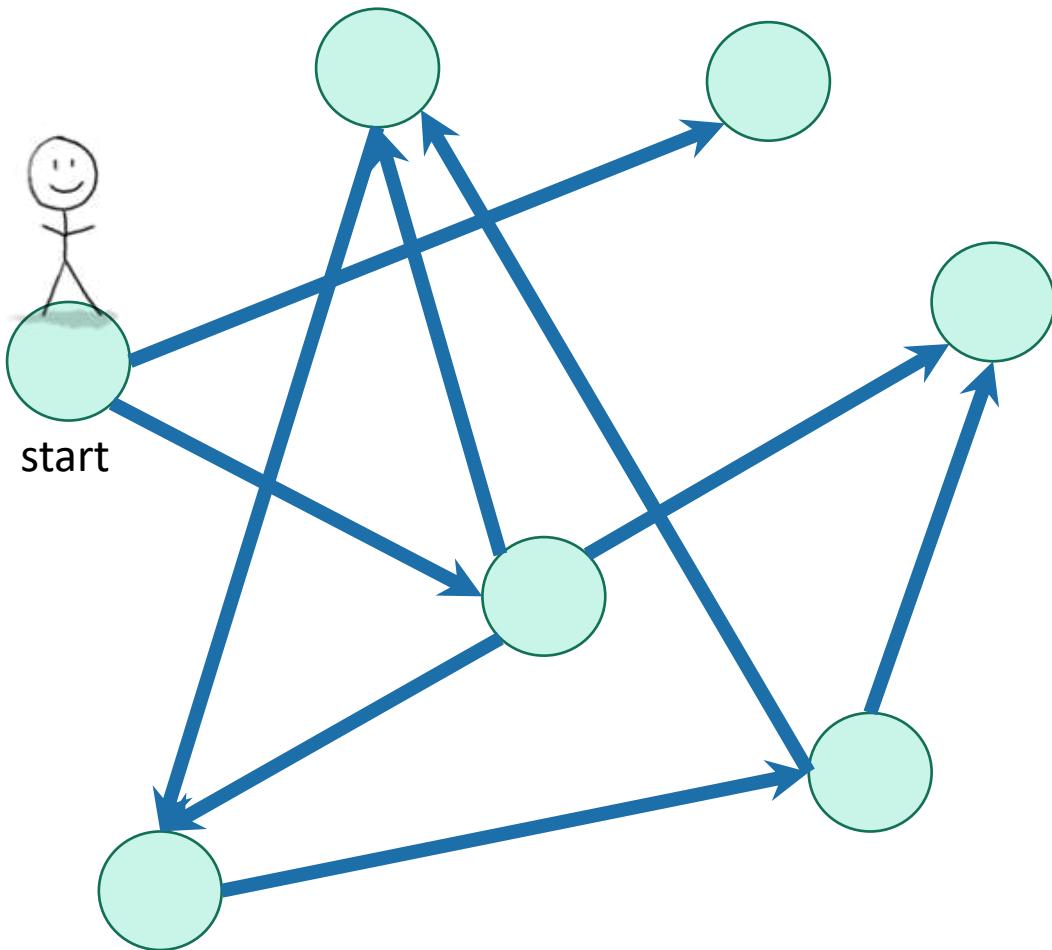
Recall: DFS

It's how you'd explore a labyrinth with chalk and a piece of string.



Depth First Search

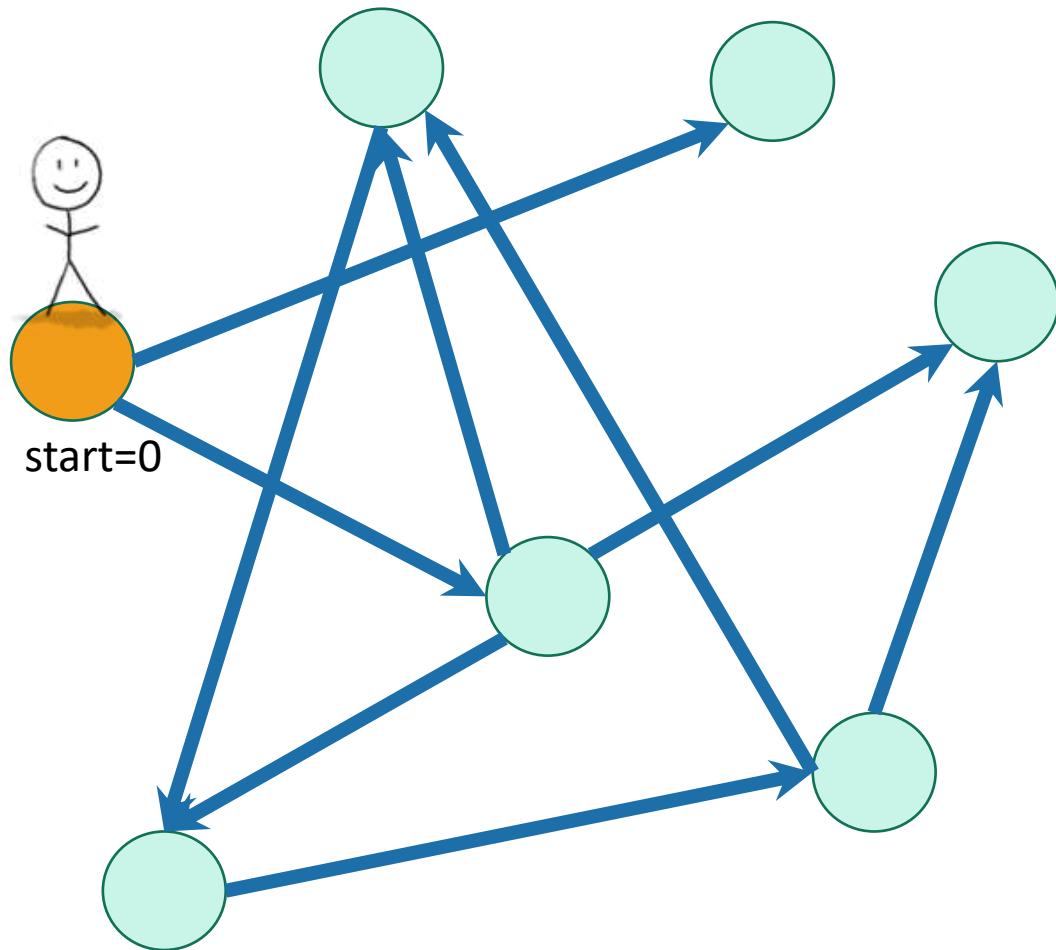
Exploring a labyrinth with chalk and a piece of string



This is the same picture we had Monday, except I've directed all of the edges. Notice that there **ARE** cycles.

Depth First Search

Exploring a labyrinth with chalk and a piece of string



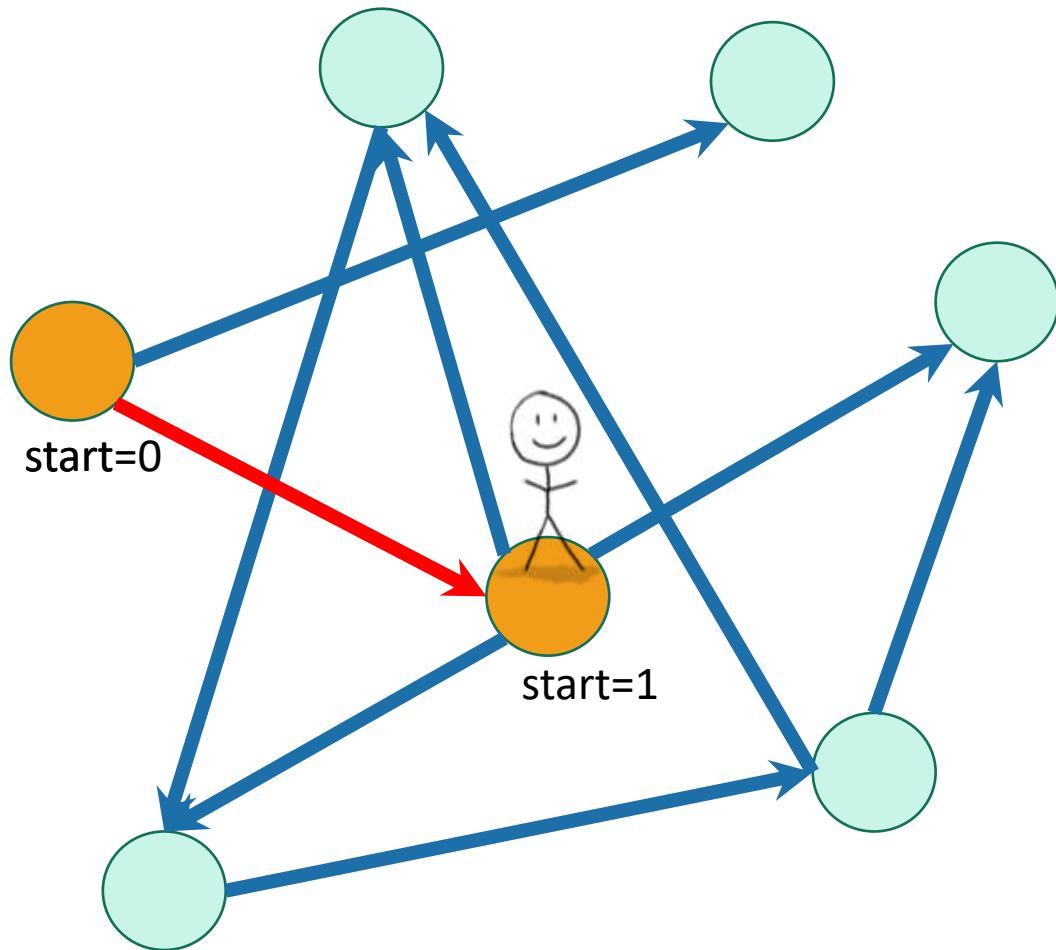
- Not been there yet
- Been there, haven't explored all the paths out.
- Been there, have explored all the paths out.



Recall we also keep track of **start** and **finish** times for every node.

Depth First Search

Exploring a labyrinth with chalk and a piece of string



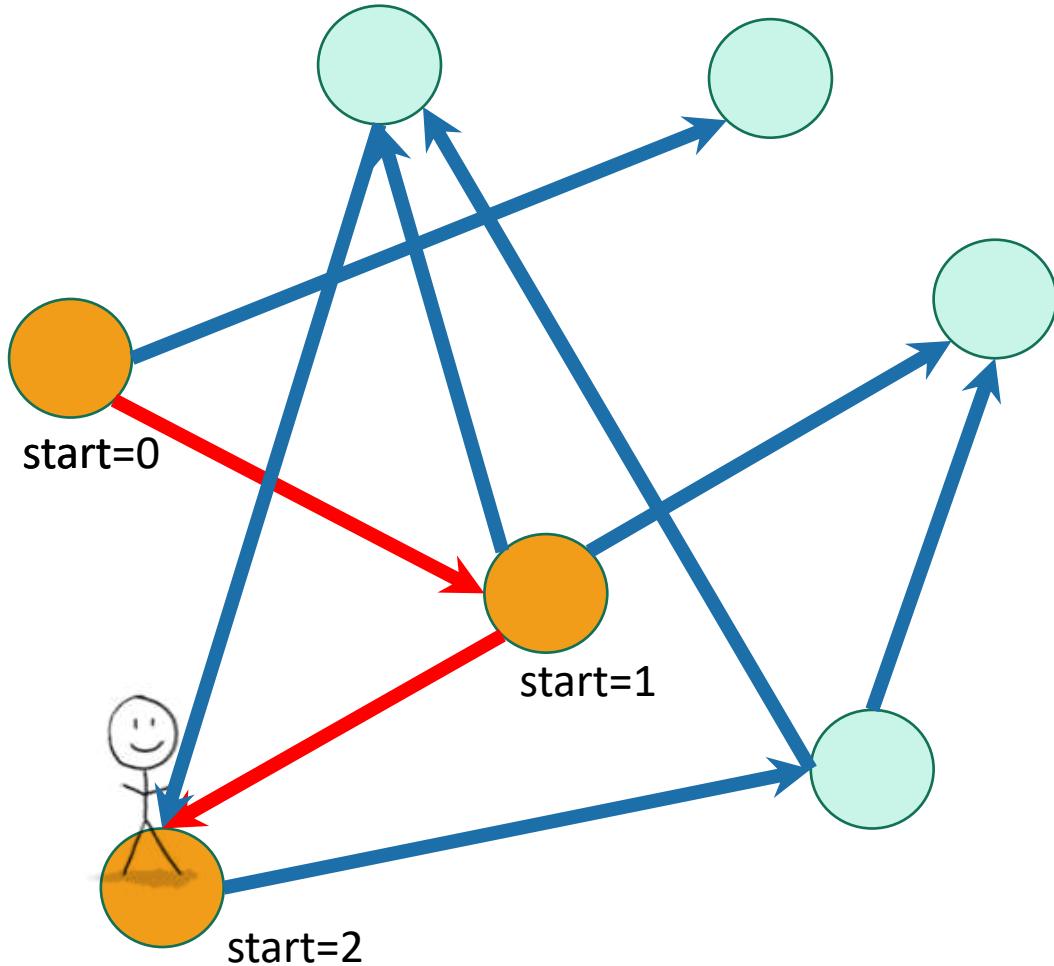
- Not been there yet
- Been there, haven't explored all the paths out.
- Been there, have explored all the paths out.



Recall we also keep track of **start** and **finish** times for every node.

Depth First Search

Exploring a labyrinth with chalk and a piece of string



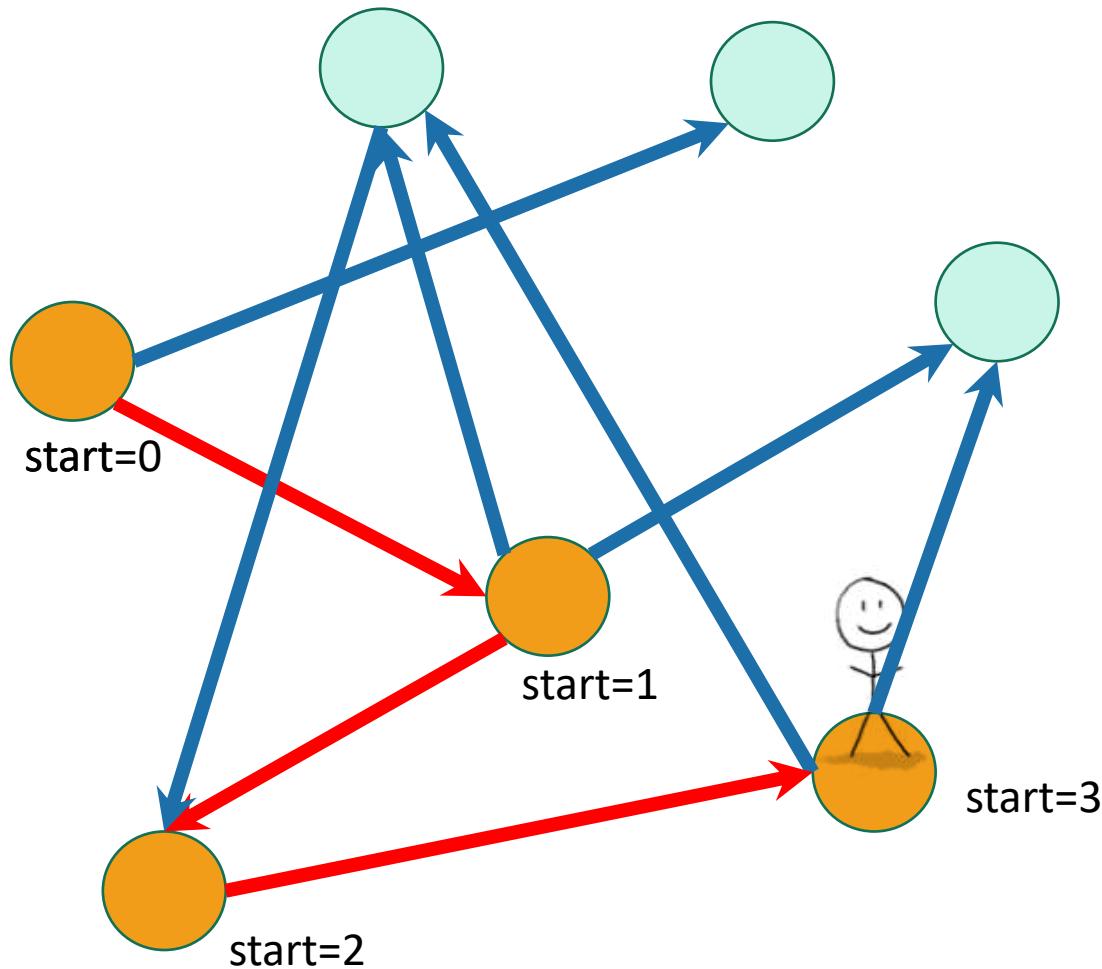
- (Light Green Circle) Not been there yet
- (Orange Circle) Been there, haven't explored all the paths out.
- (Dark Green Circle) Been there, have explored all the paths out.



Recall we also keep track of **start** and **finish** times for every node.

Depth First Search

Exploring a labyrinth with chalk and a piece of string



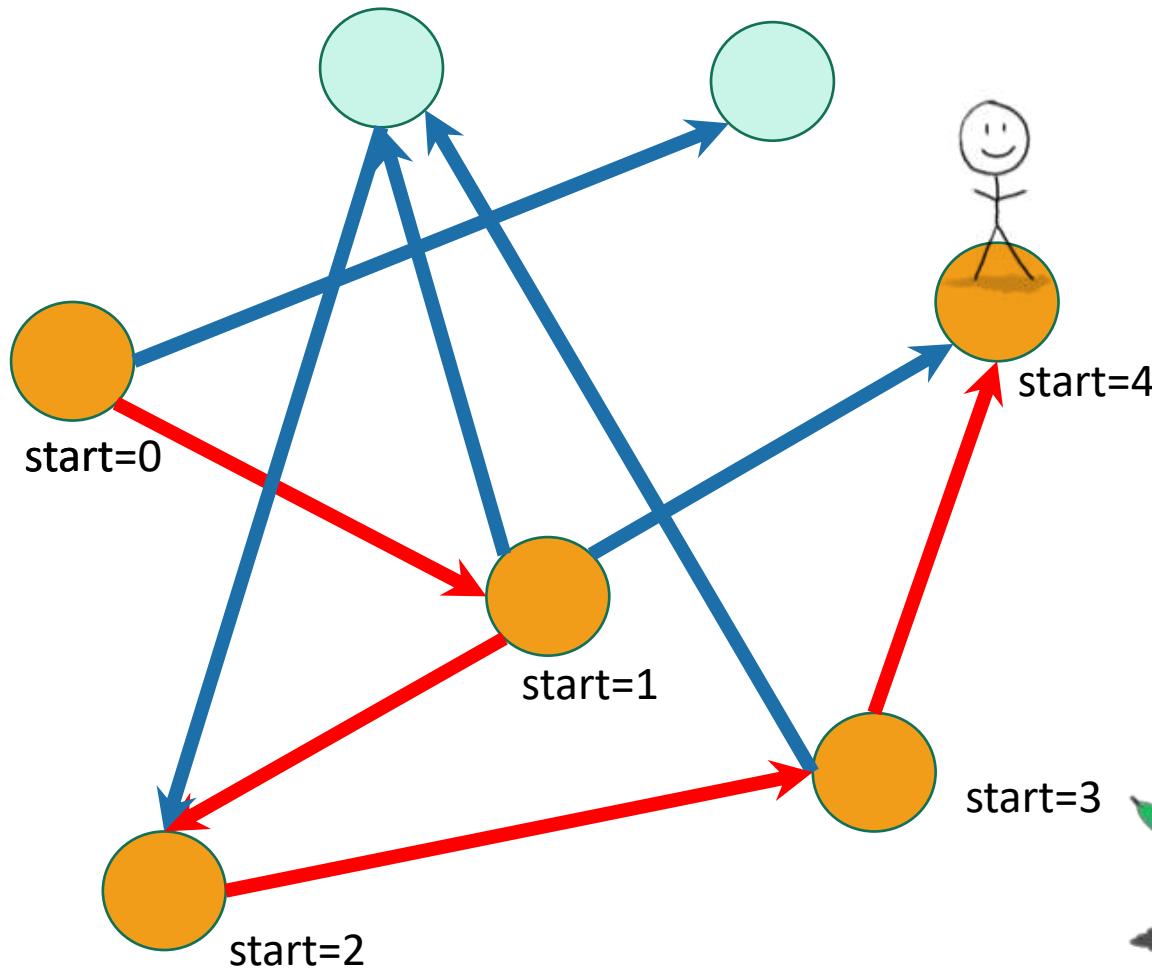
- Not been there yet
- Been there, haven't explored all the paths out.
- Been there, have explored all the paths out.



Recall we also keep track of **start** and **finish** times for every node.

Depth First Search

Exploring a labyrinth with chalk and a piece of string



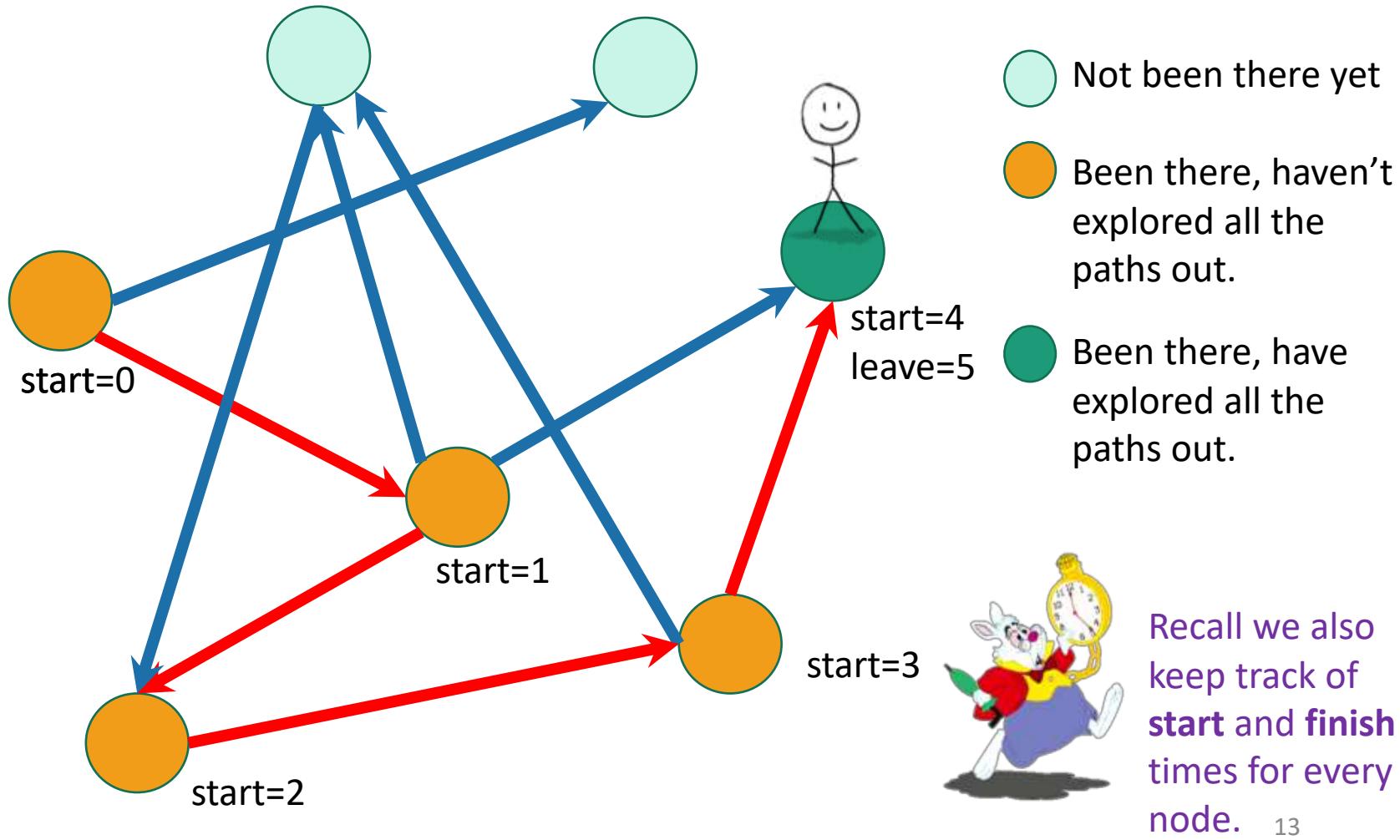
- (Light Green Circle) Not been there yet
- (Orange Circle) Been there, haven't explored all the paths out.
- (Dark Green Circle) Been there, have explored all the paths out.



Recall we also keep track of **start** and **finish** times for every node.

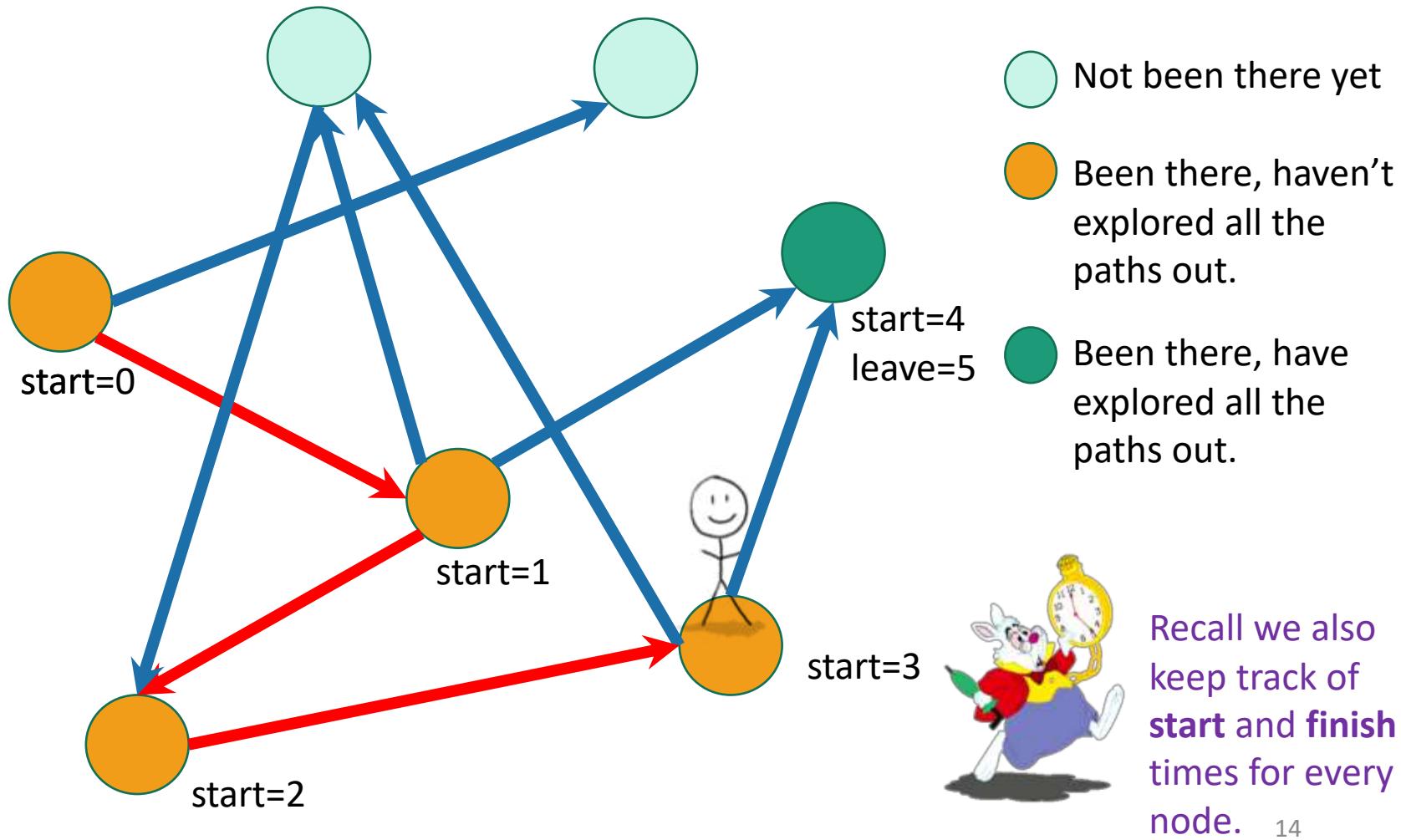
Depth First Search

Exploring a labyrinth with chalk and a piece of string



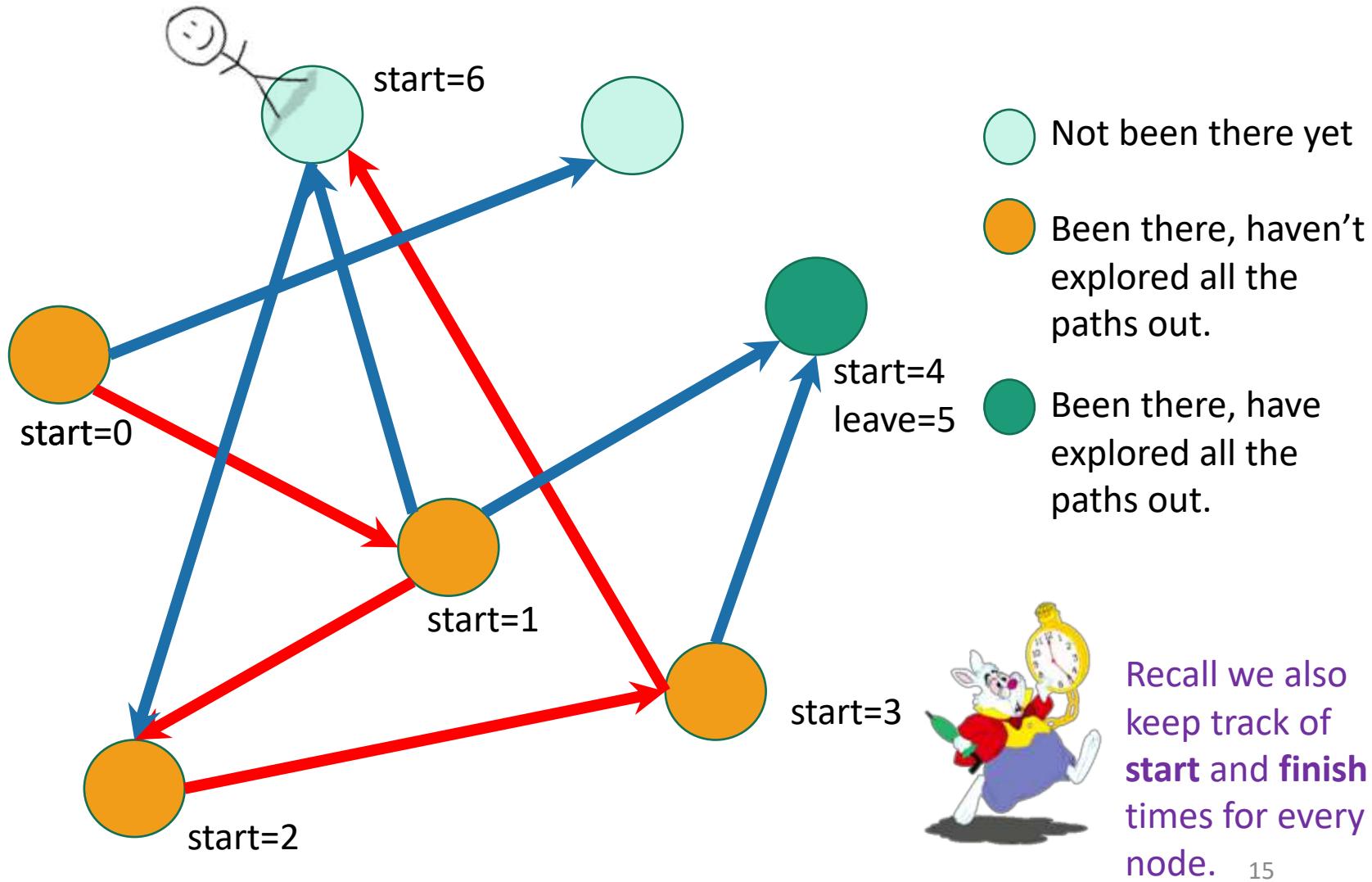
Depth First Search

Exploring a labyrinth with chalk and a piece of string



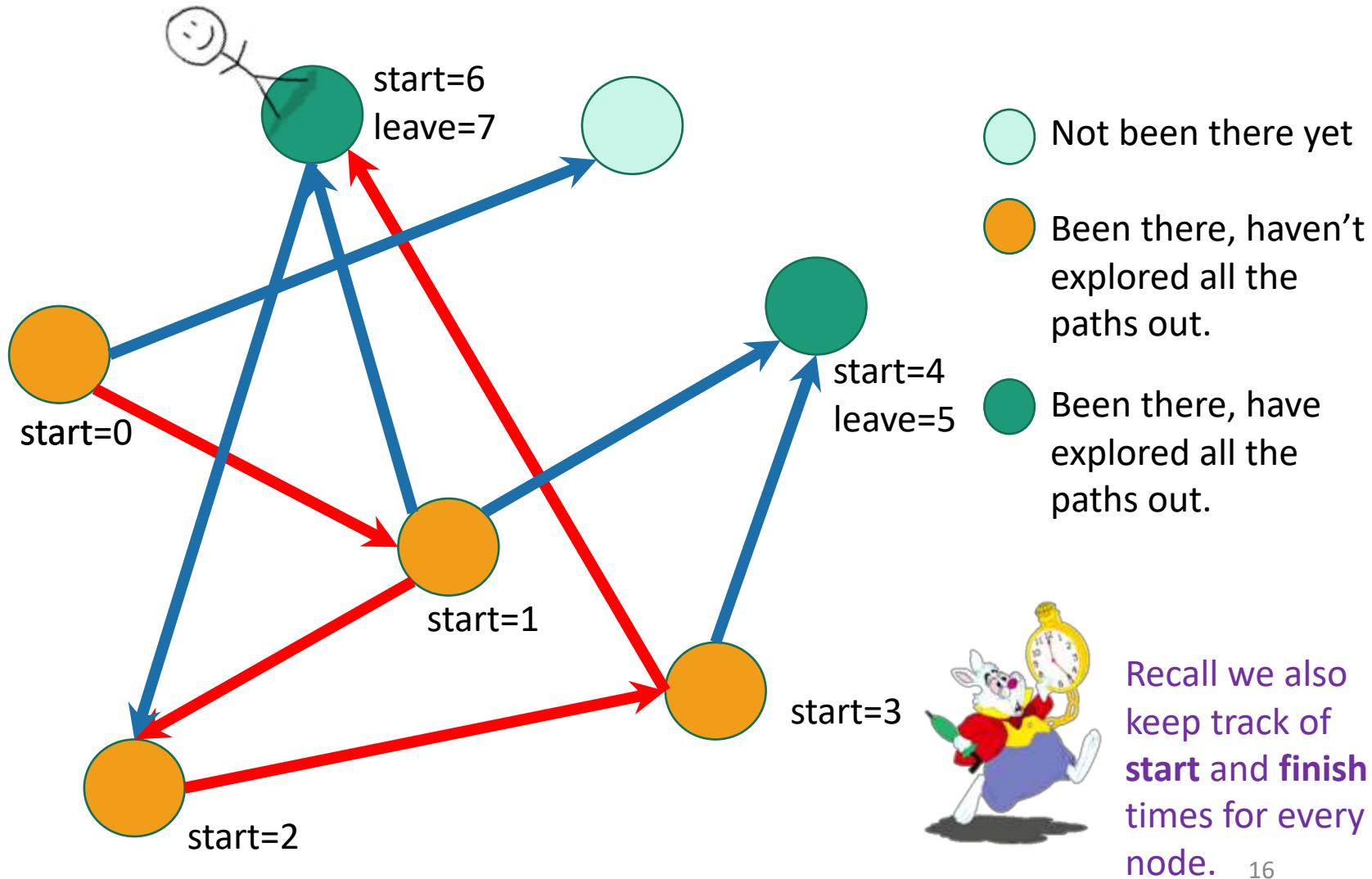
Depth First Search

Exploring a labyrinth with chalk and a piece of string



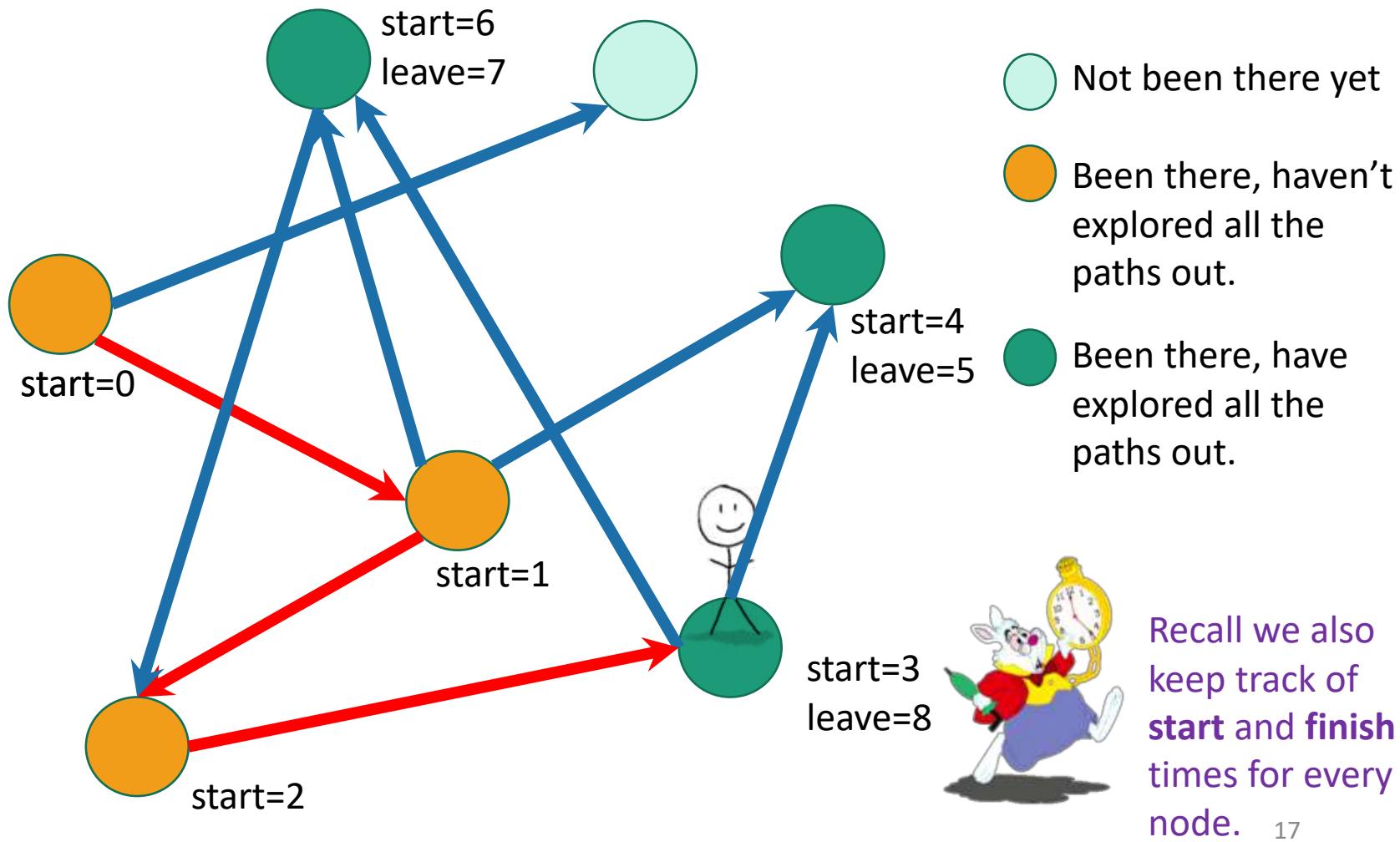
Depth First Search

Exploring a labyrinth with chalk and a piece of string



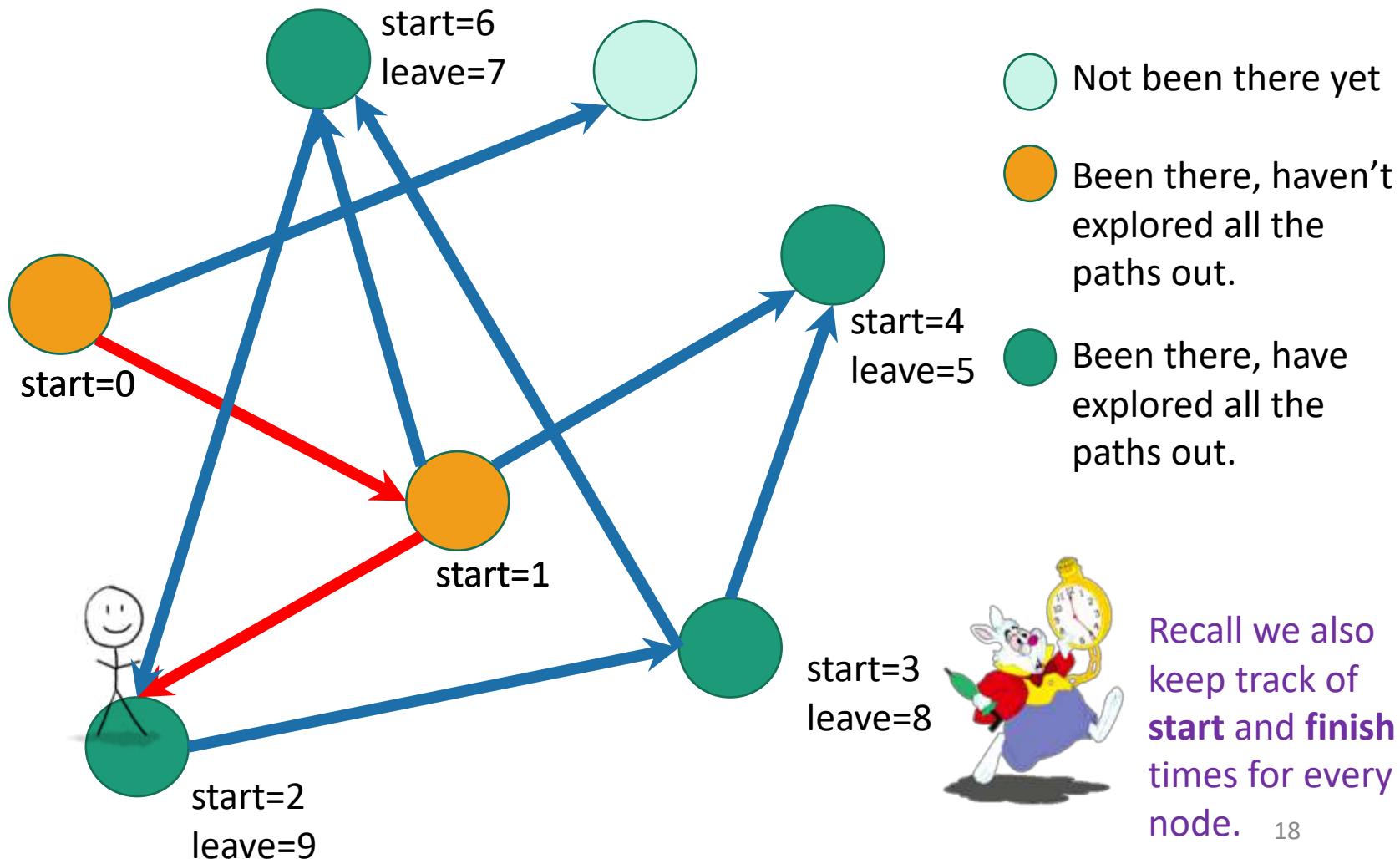
Depth First Search

Exploring a labyrinth with chalk and a piece of string



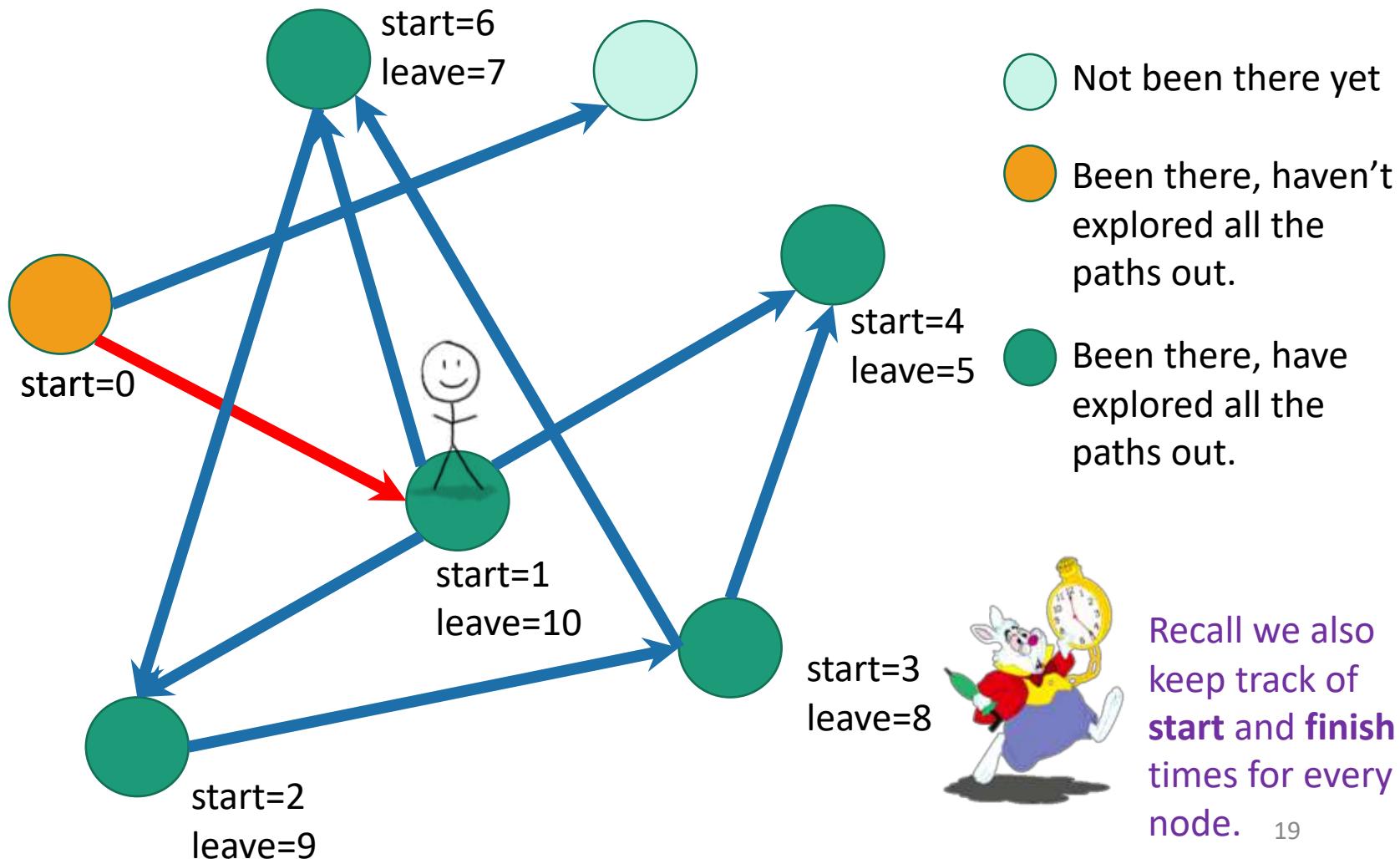
Depth First Search

Exploring a labyrinth with chalk and a piece of string



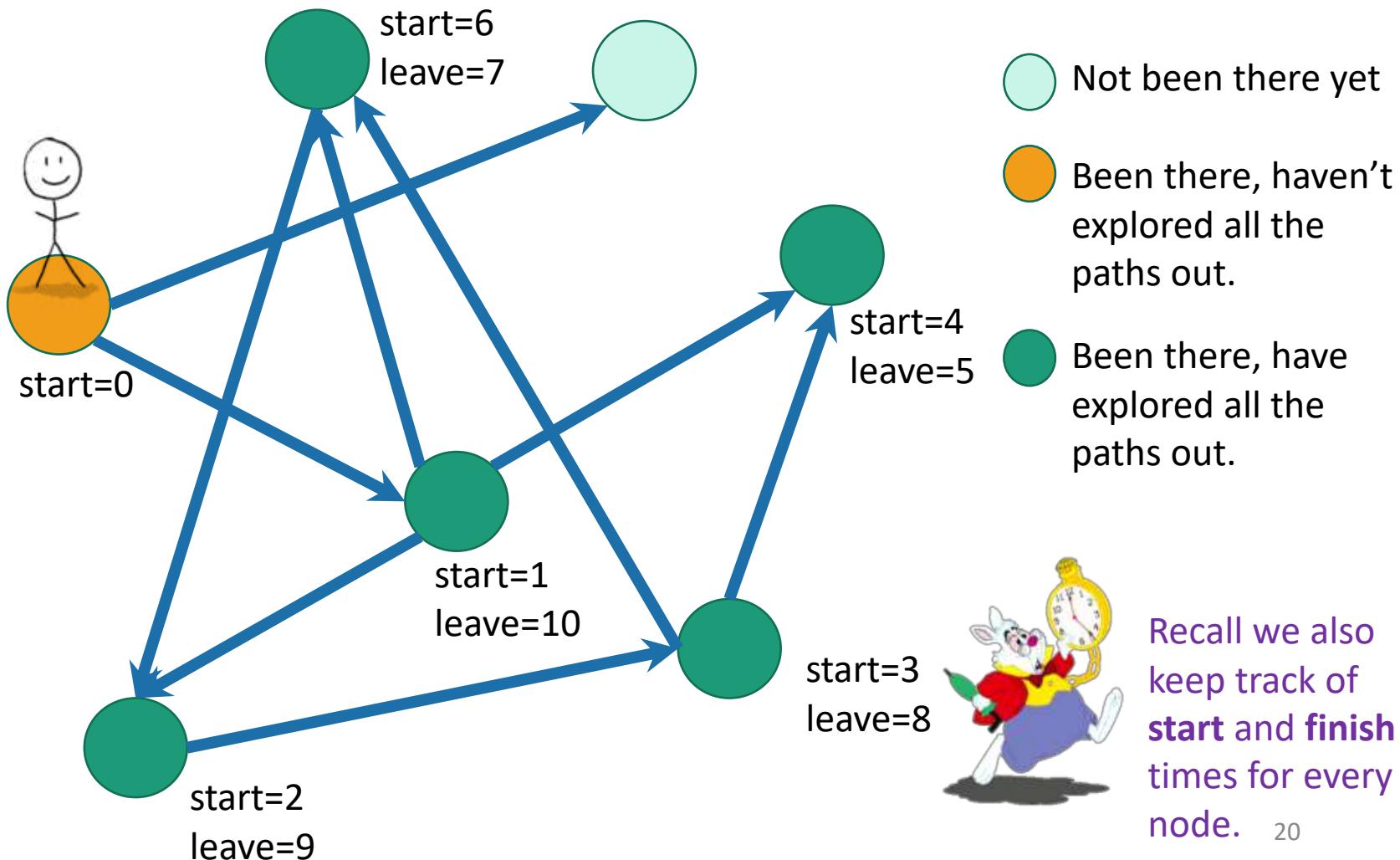
Depth First Search

Exploring a labyrinth with chalk and a piece of string



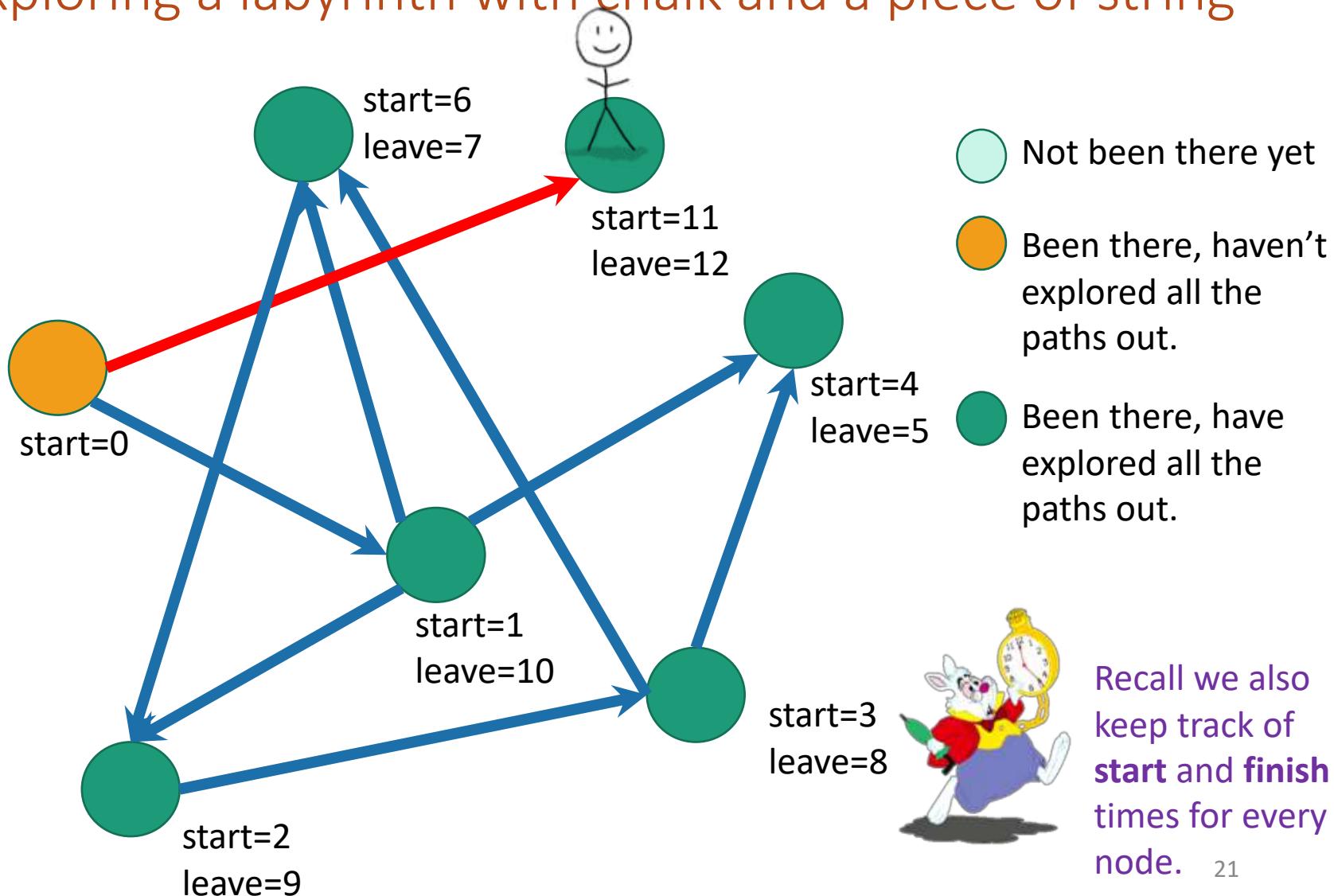
Depth First Search

Exploring a labyrinth with chalk and a piece of string



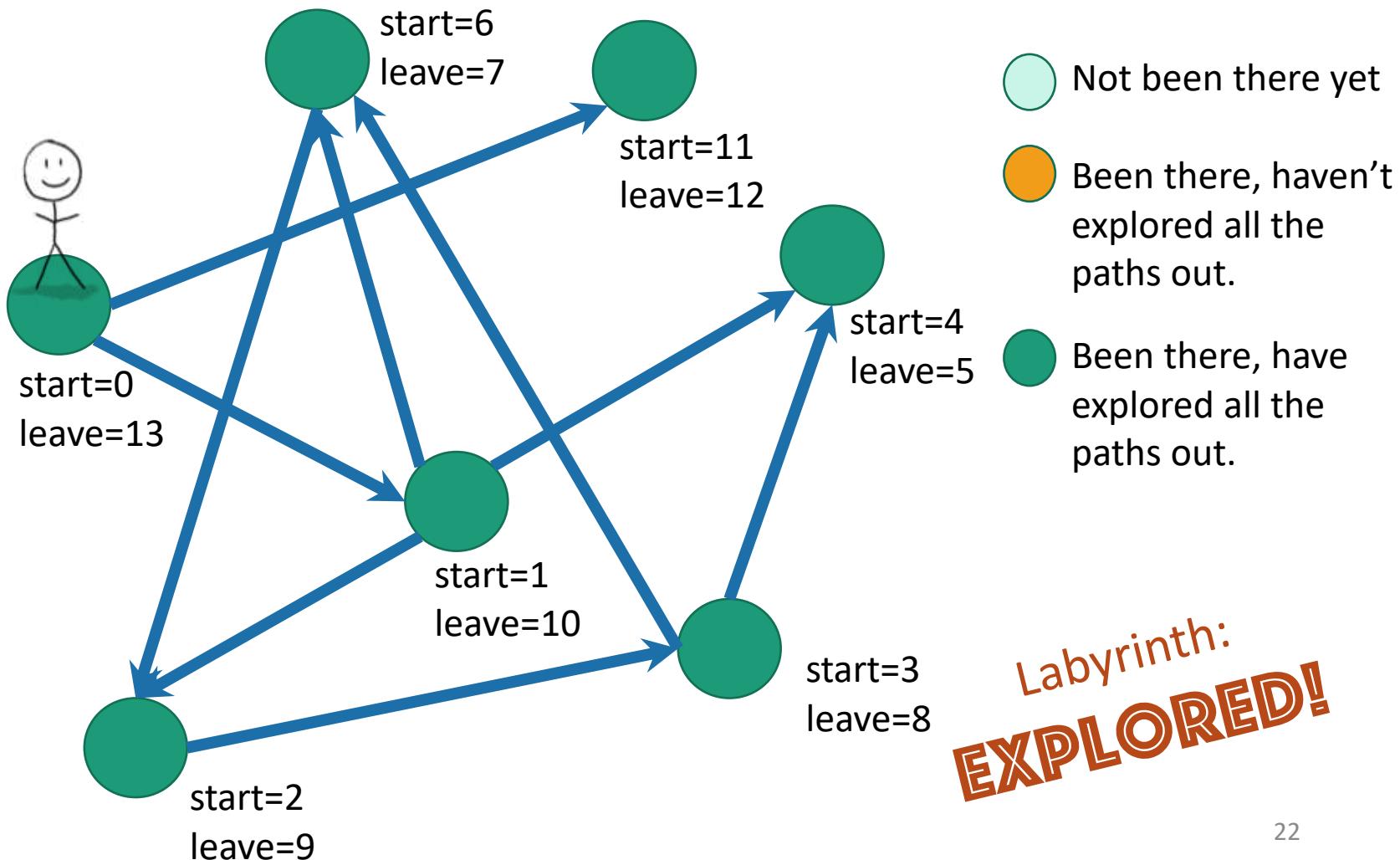
Depth First Search

Exploring a labyrinth with chalk and a piece of string



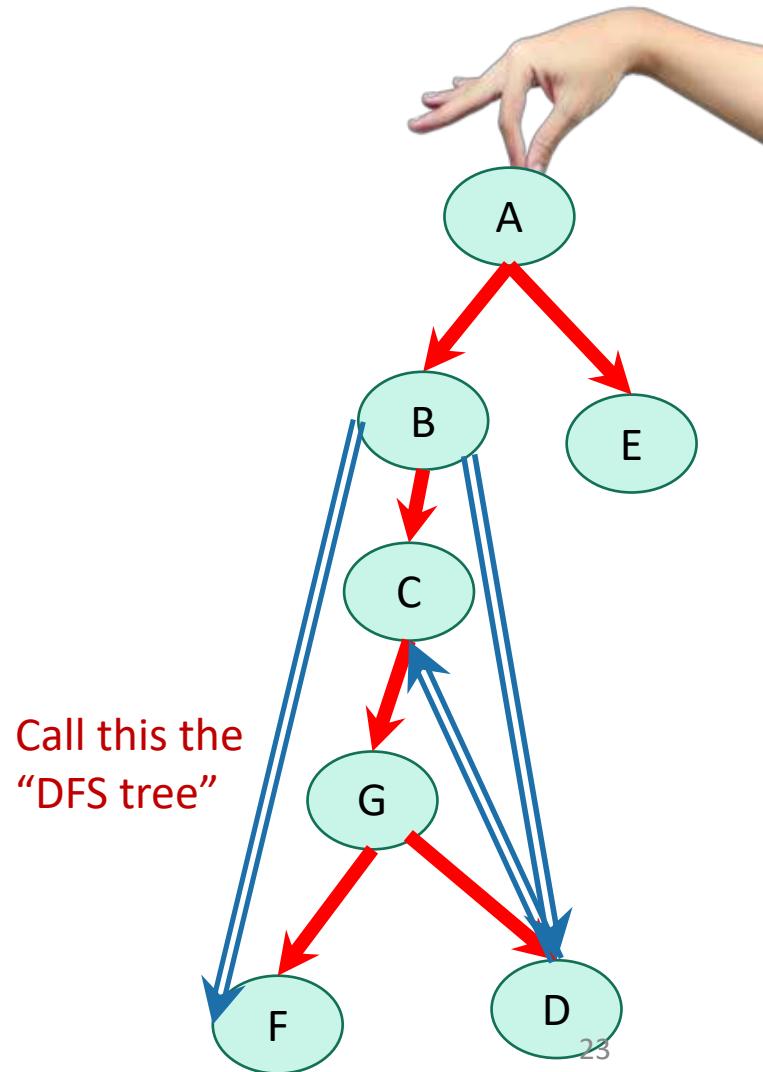
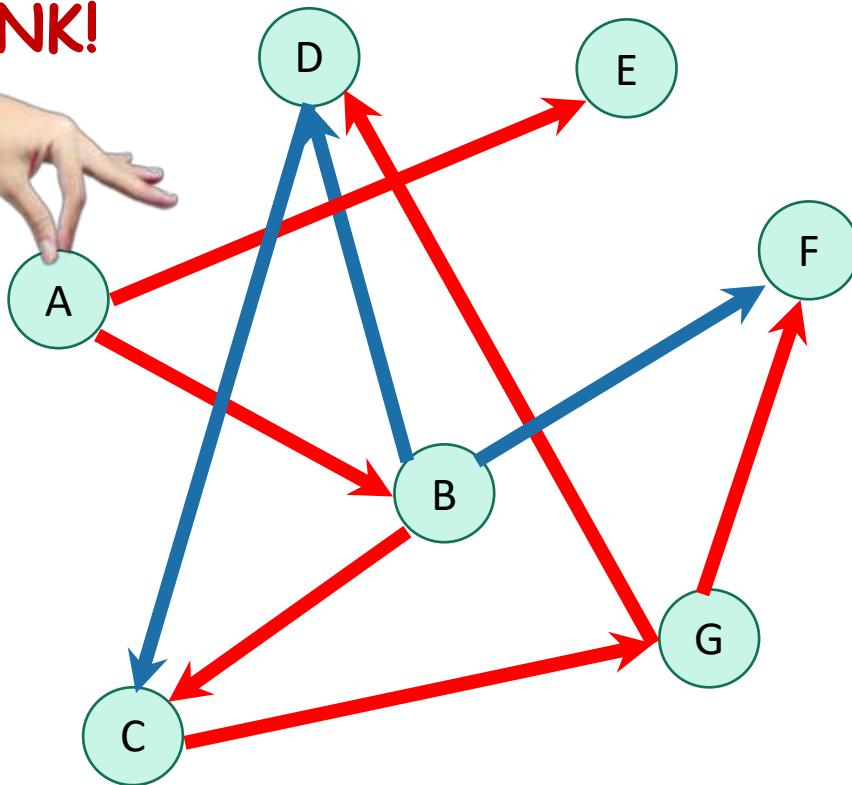
Depth First Search

Exploring a labyrinth with chalk and a piece of string



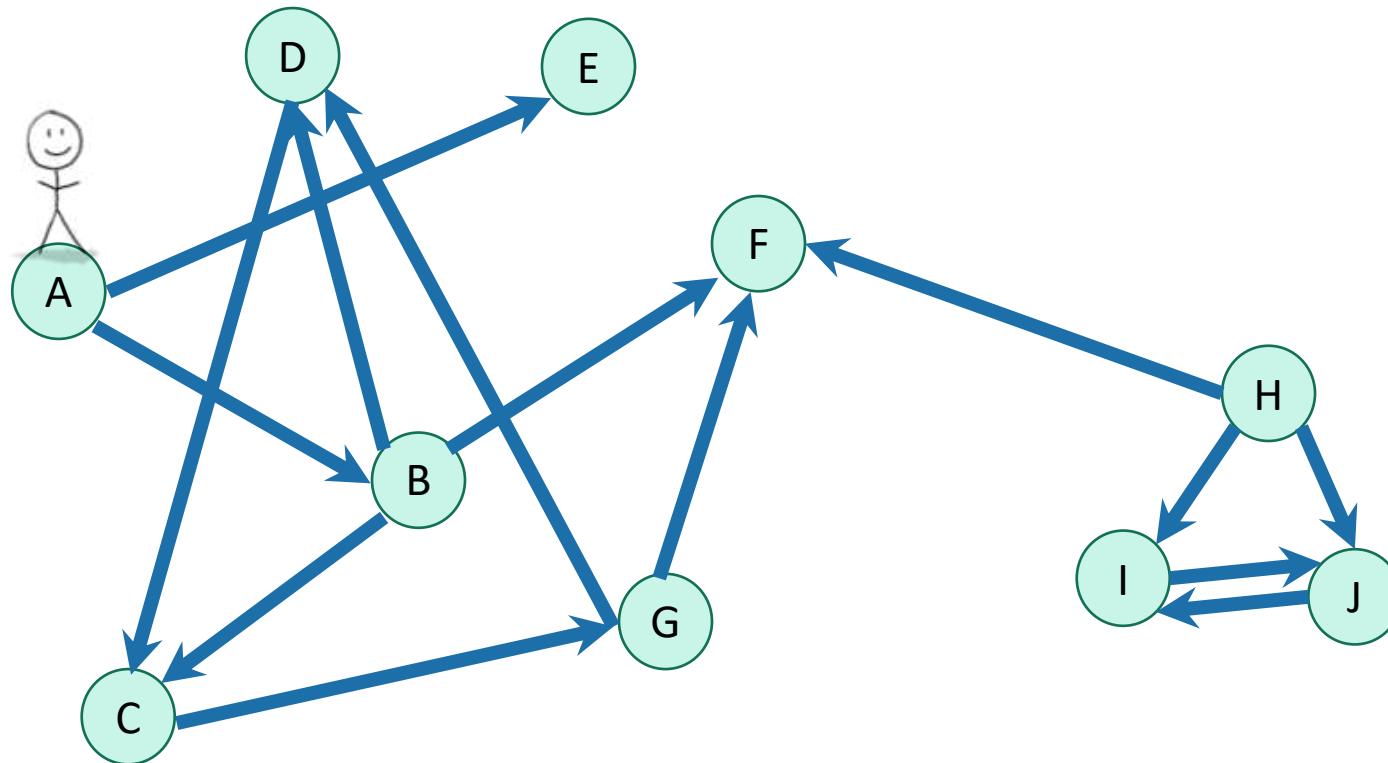
Depth first search

implicitly creates a tree on everything you can reach



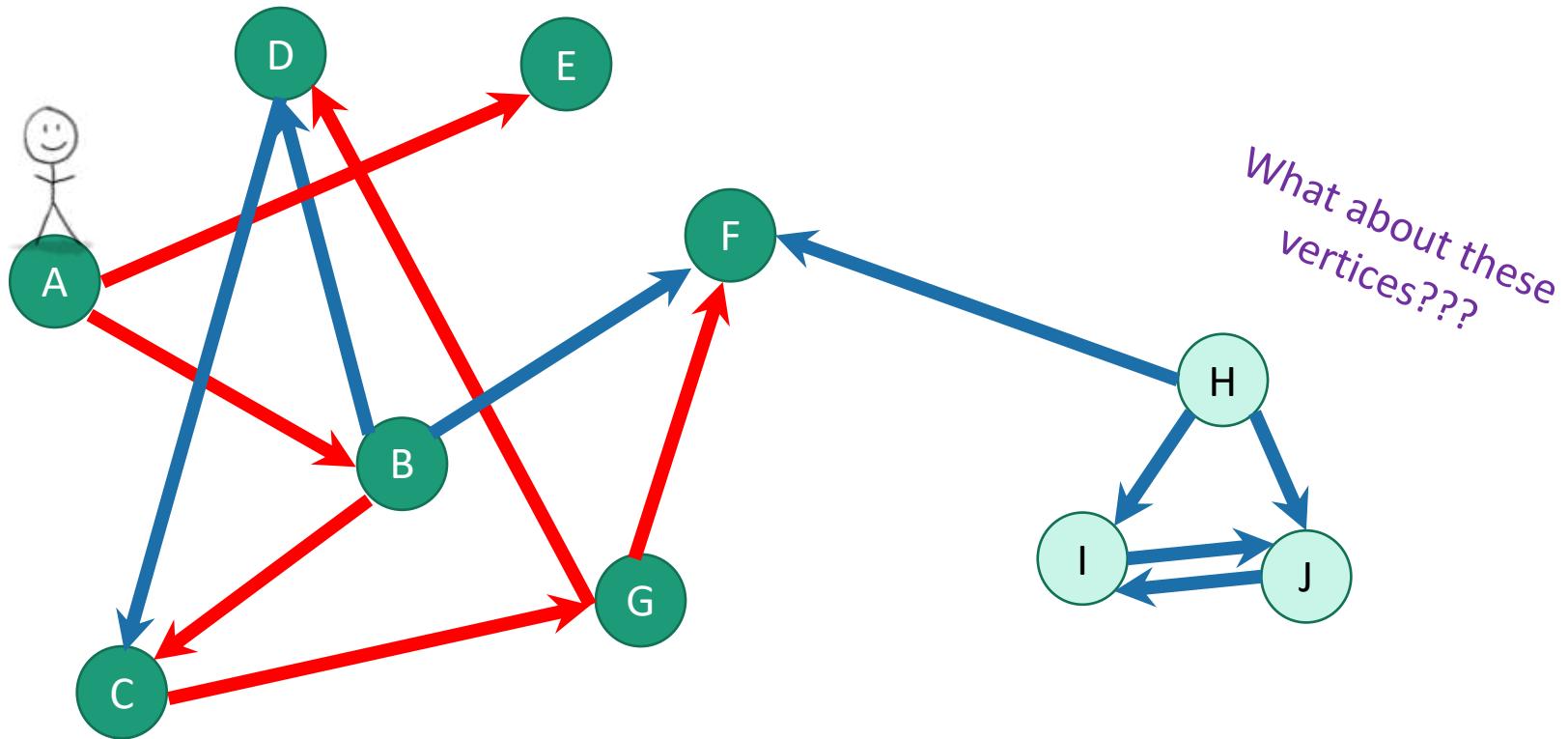
When you can't reach everything

- Run DFS repeatedly to get a **depth-first forest**



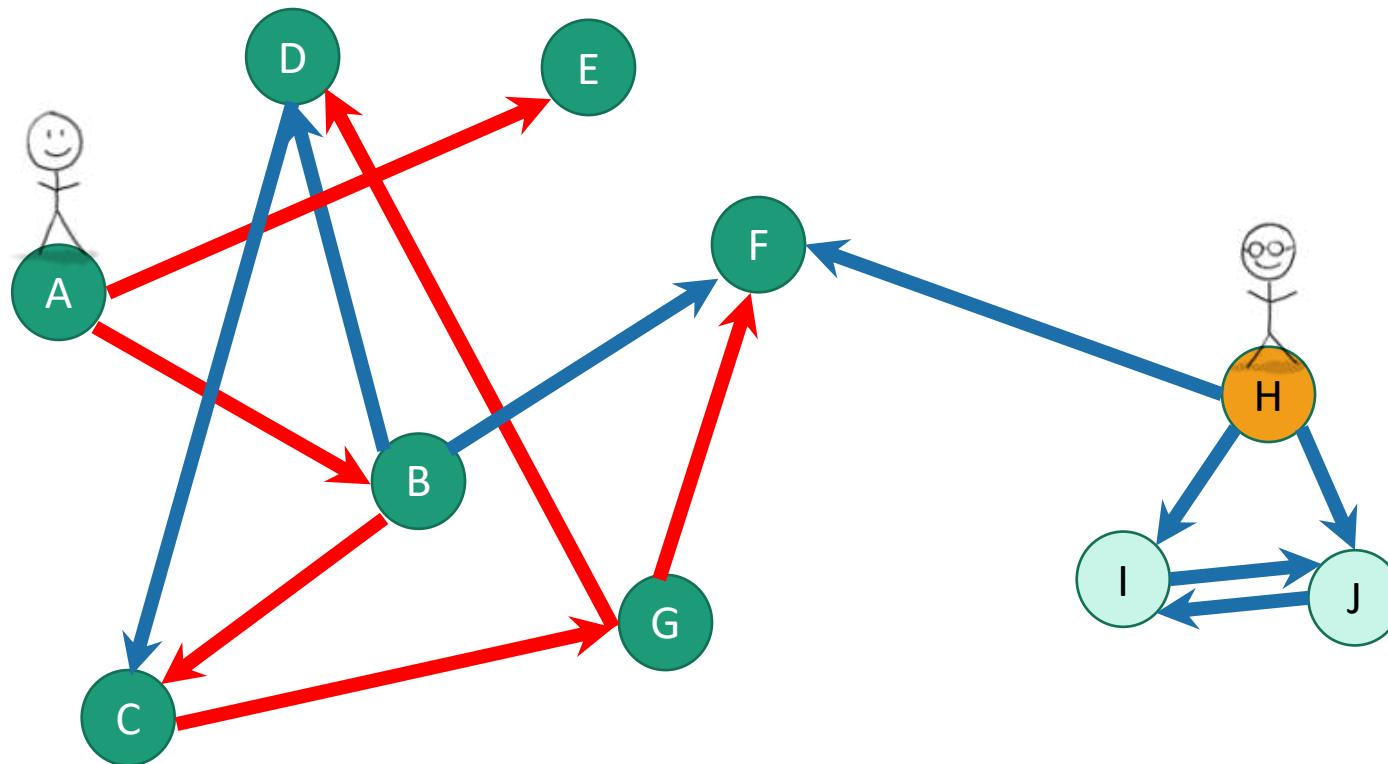
When you can't reach everything

- Run DFS repeatedly to get a **depth-first forest**



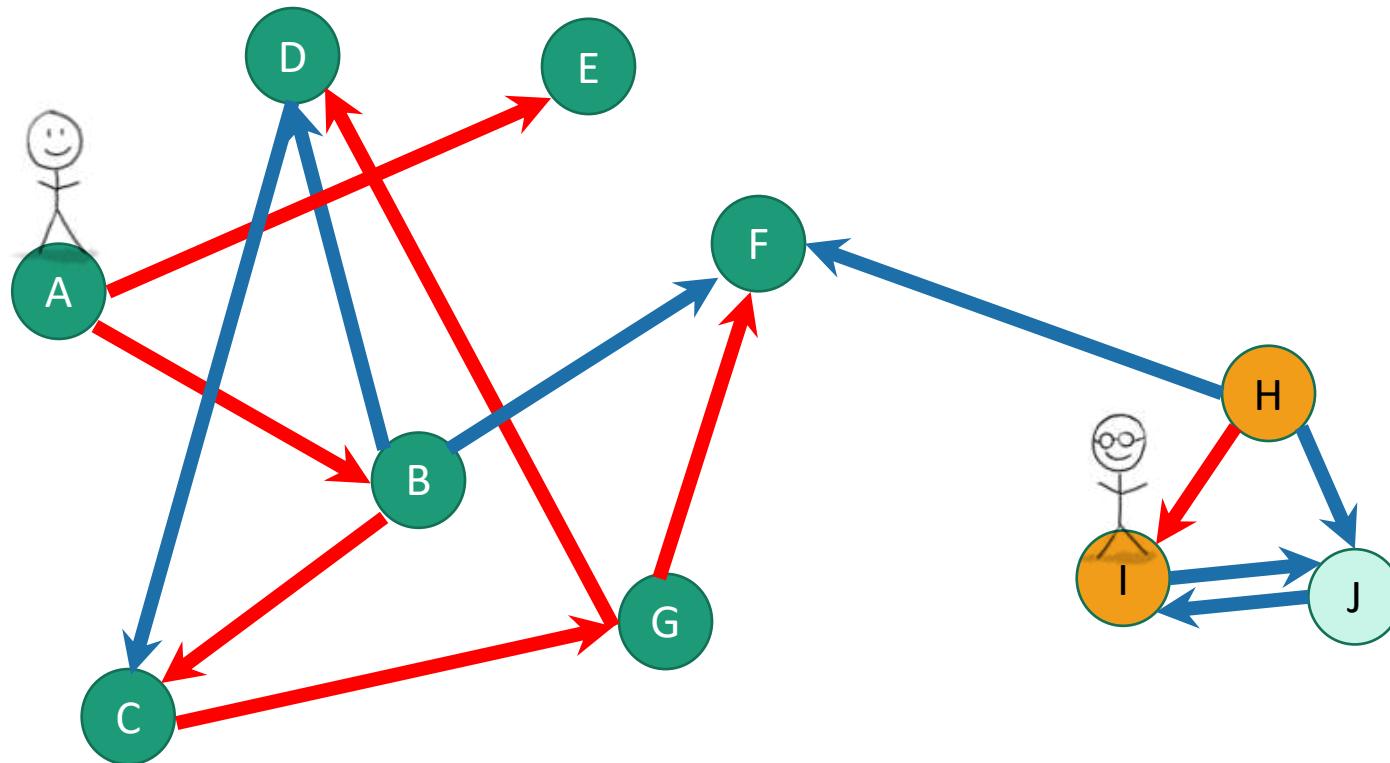
When you can't reach everything

- Run DFS repeatedly to get a **depth-first forest**



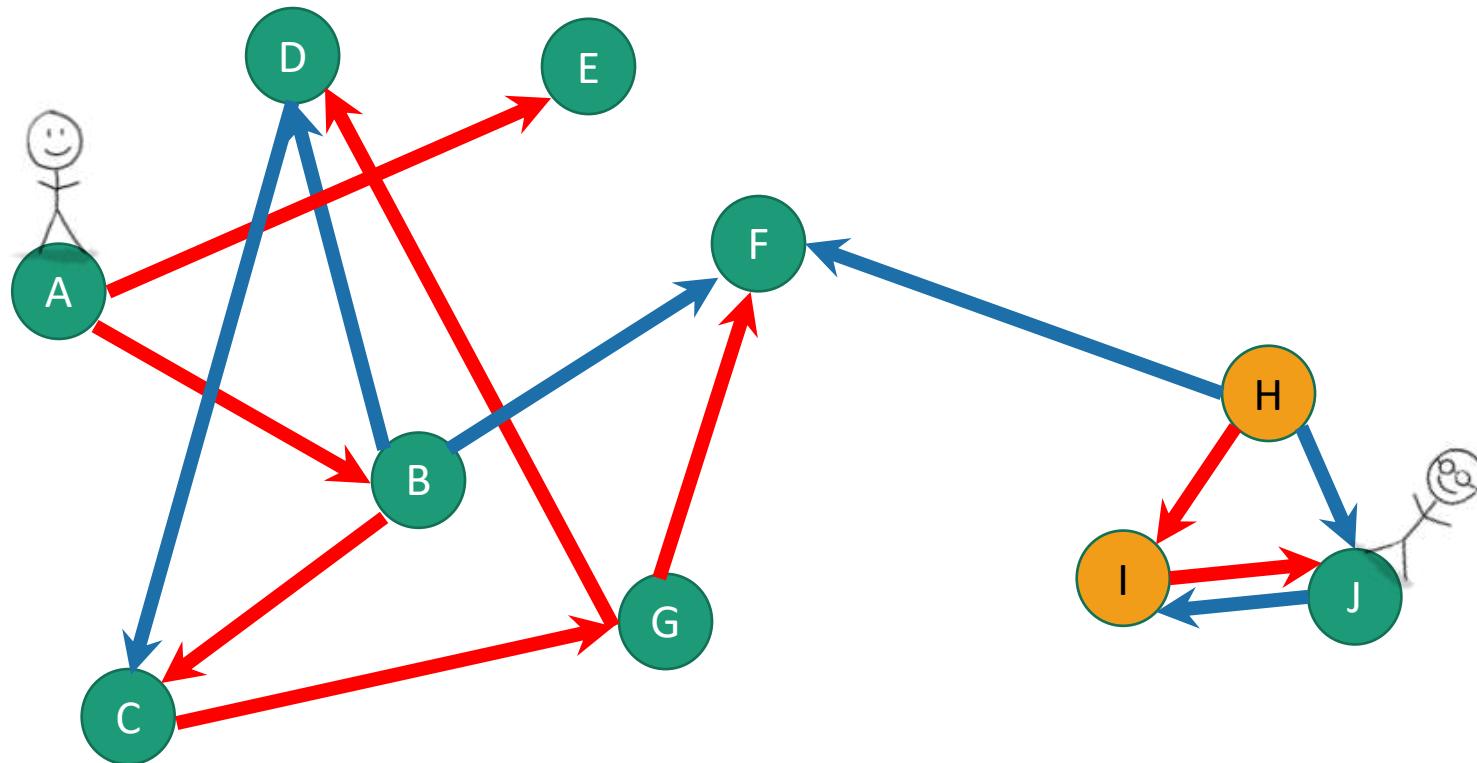
When you can't reach everything

- Run DFS repeatedly to get a **depth-first forest**



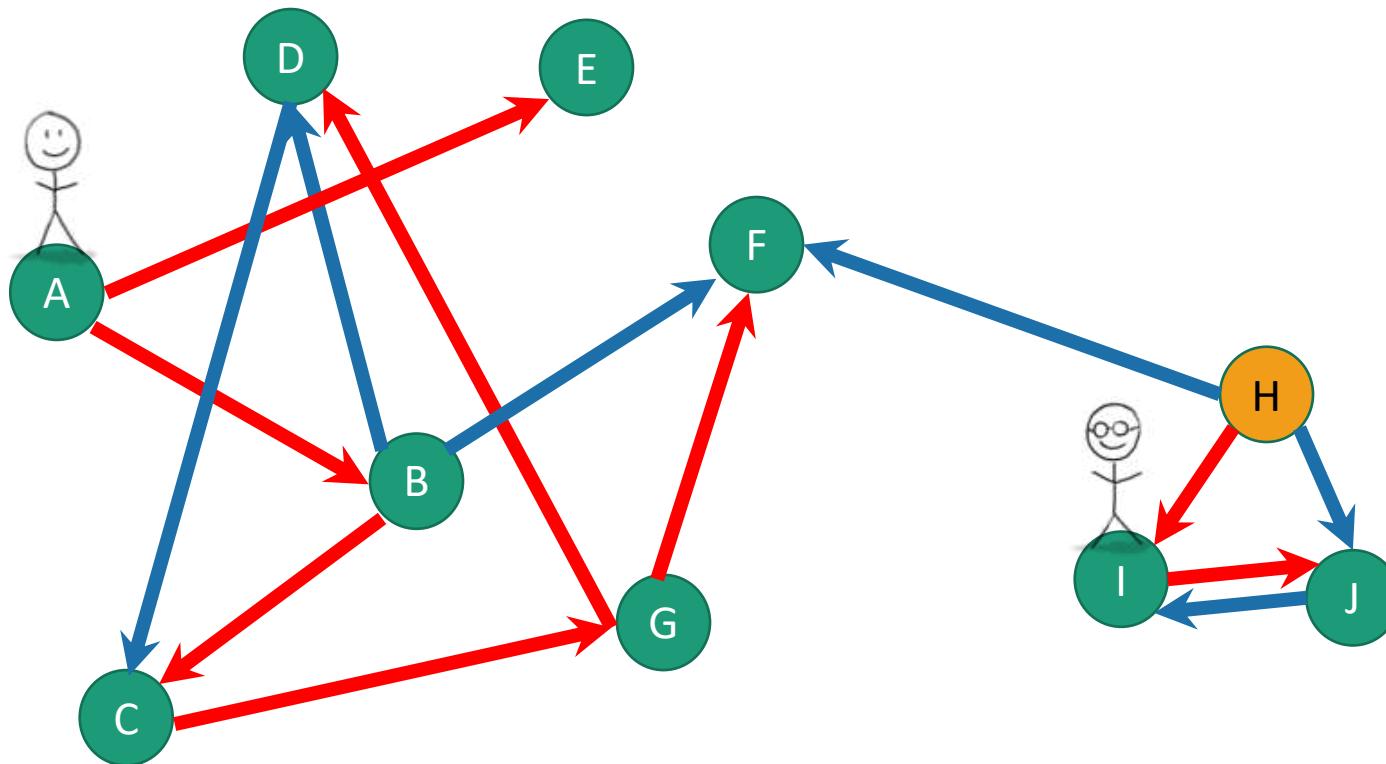
When you can't reach everything

- Run DFS repeatedly to get a **depth-first forest**



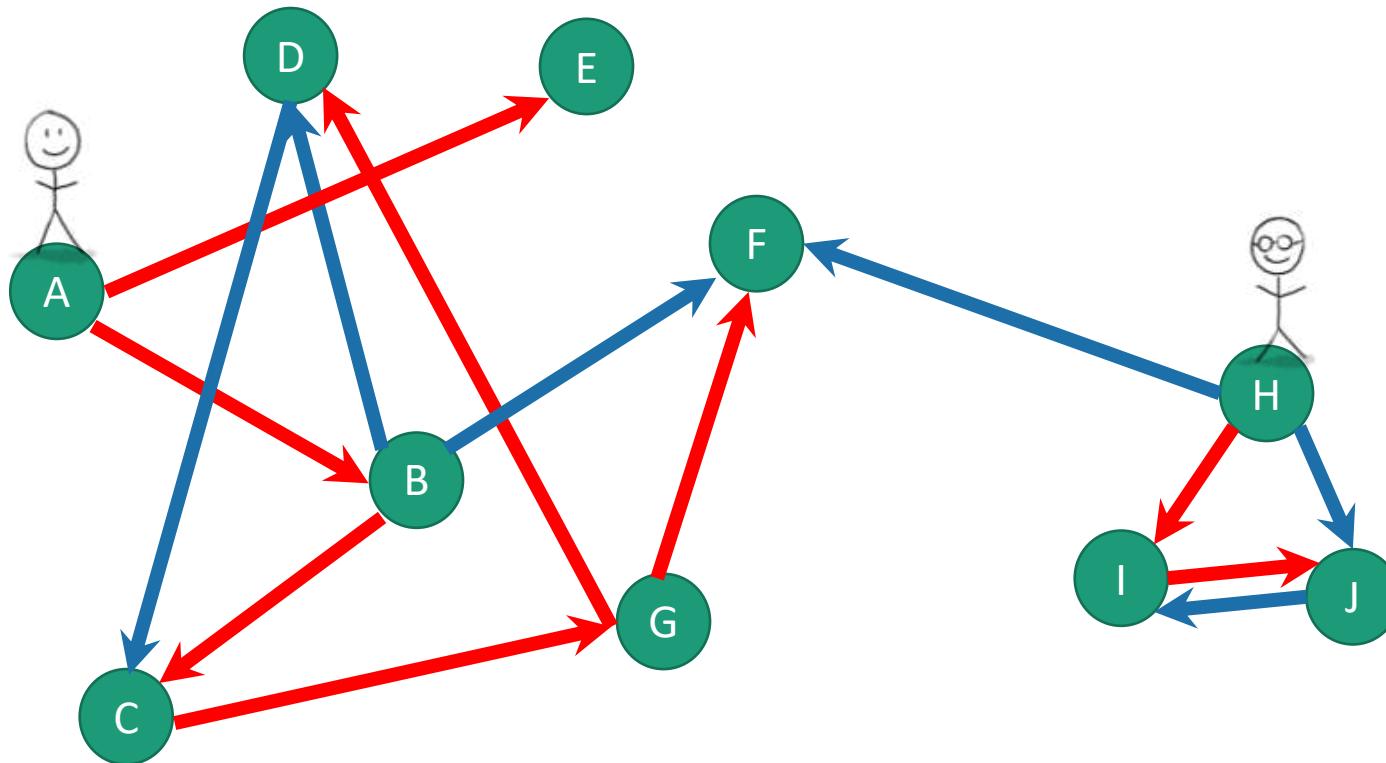
When you can't reach everything

- Run DFS repeatedly to get a **depth-first forest**



When you can't reach everything

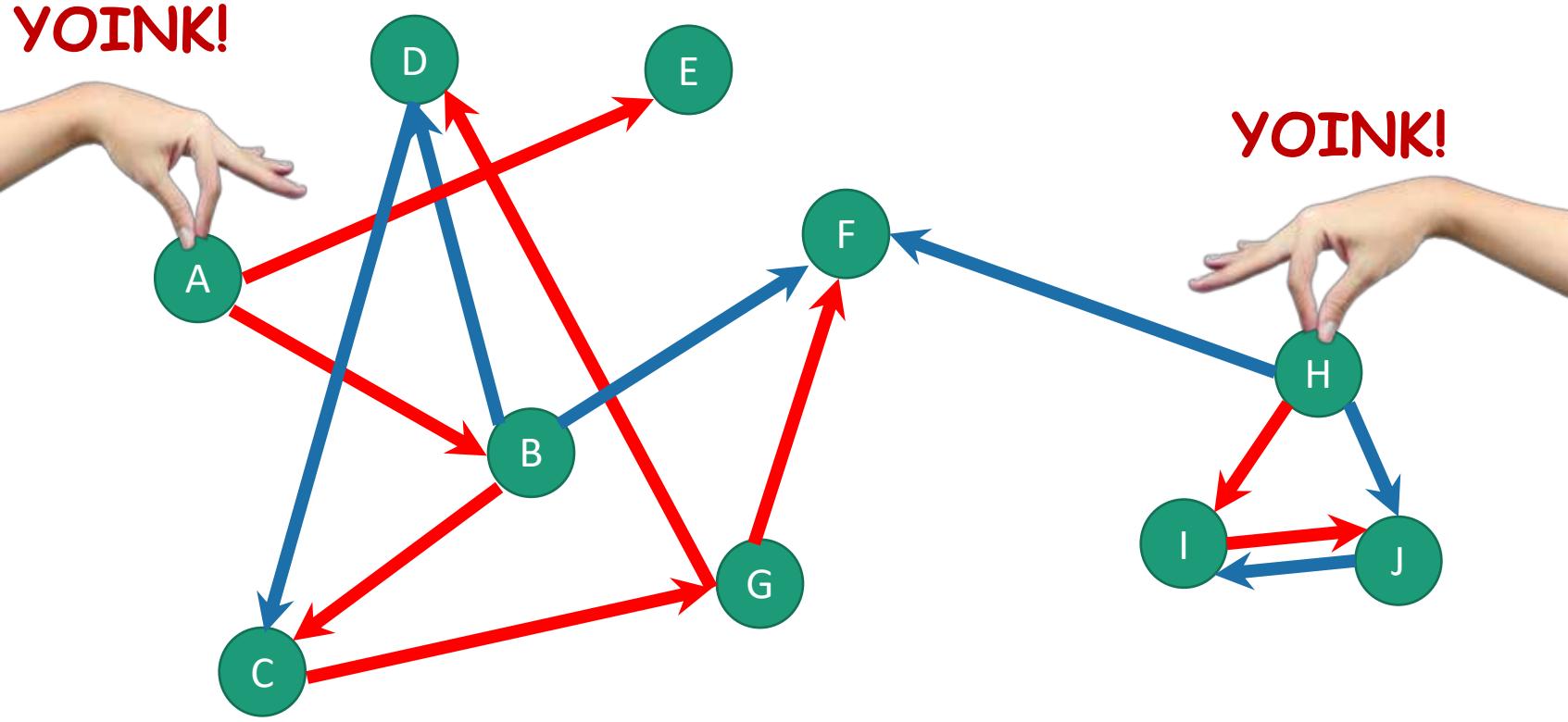
- Run DFS repeatedly to get a **depth-first forest**



When you can't reach everything

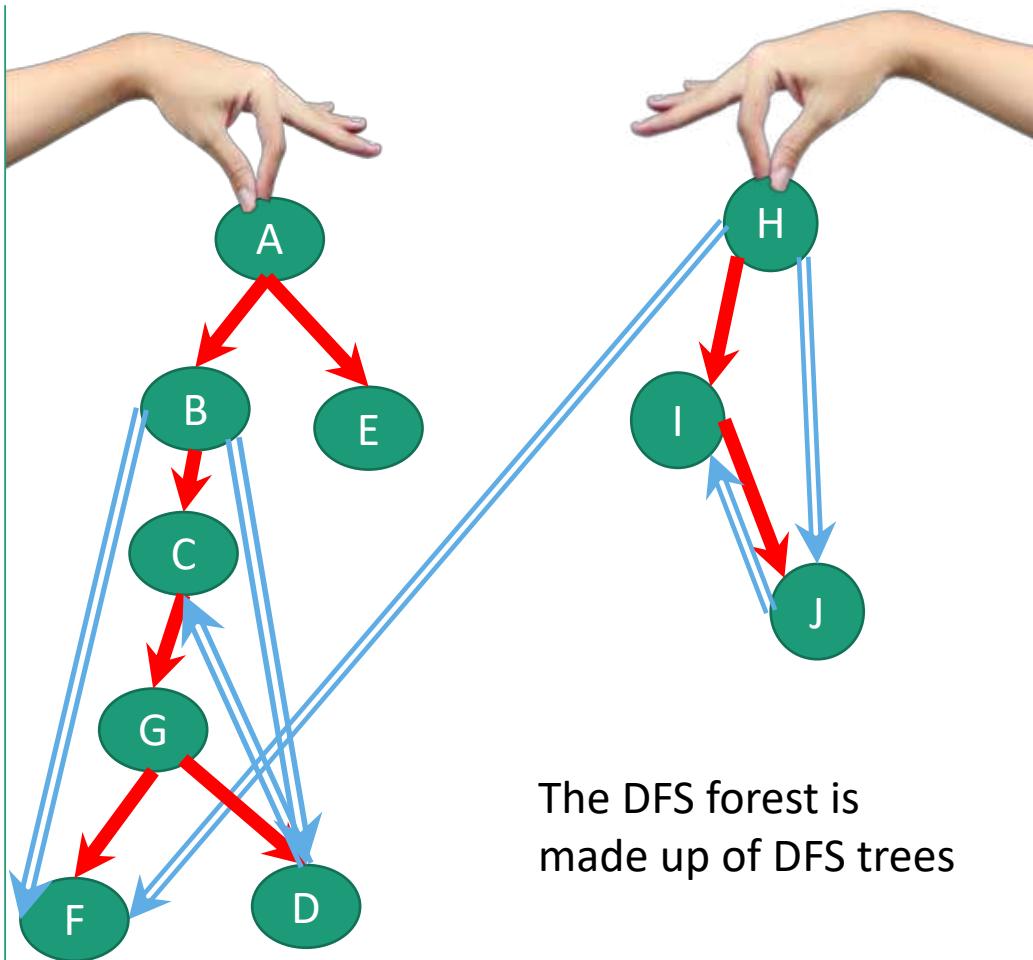
- Run DFS repeatedly to get a **depth-first forest**

YOINK!



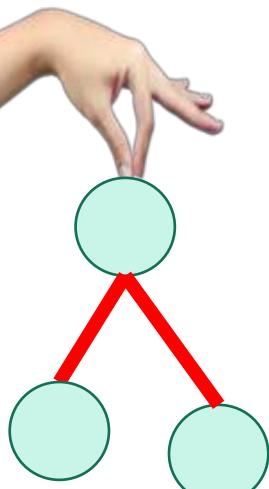
When you can't reach everything

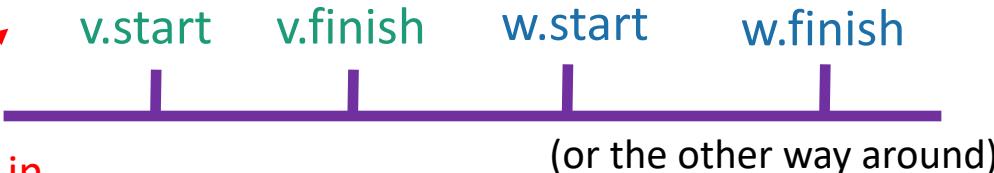
- Run DFS repeatedly to get a **depth-first forest**



Recall: the parentheses theorem

(Works the same with DFS forests)

- If v is a descendent of w in this tree:
- 
- w.start v.start v.finish w.finish
timeline
- If w is a descendent of v in this tree:
- v.start w.start w.finish v.finish
- If neither are descendants of each other:



If v and w are in different trees, it's always this last one.

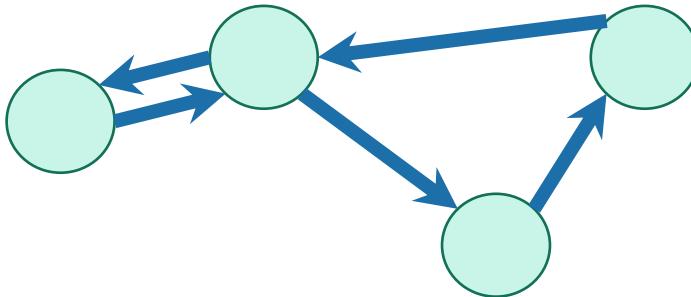
DFS tree

Enough of review

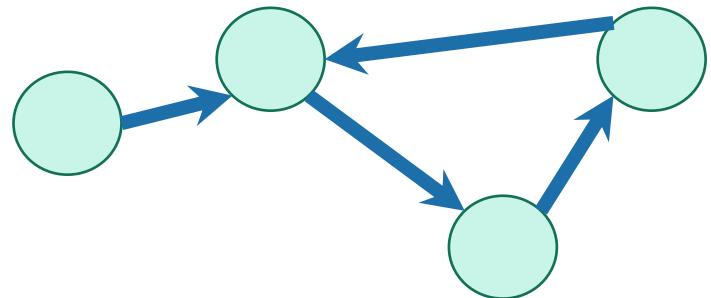
Strongly connected components

Strongly connected components

- A directed graph $G = (V, E)$ is **strongly connected** if:
- for all v, w in V :
 - there is a path from v to w and
 - there is a path from w to v .



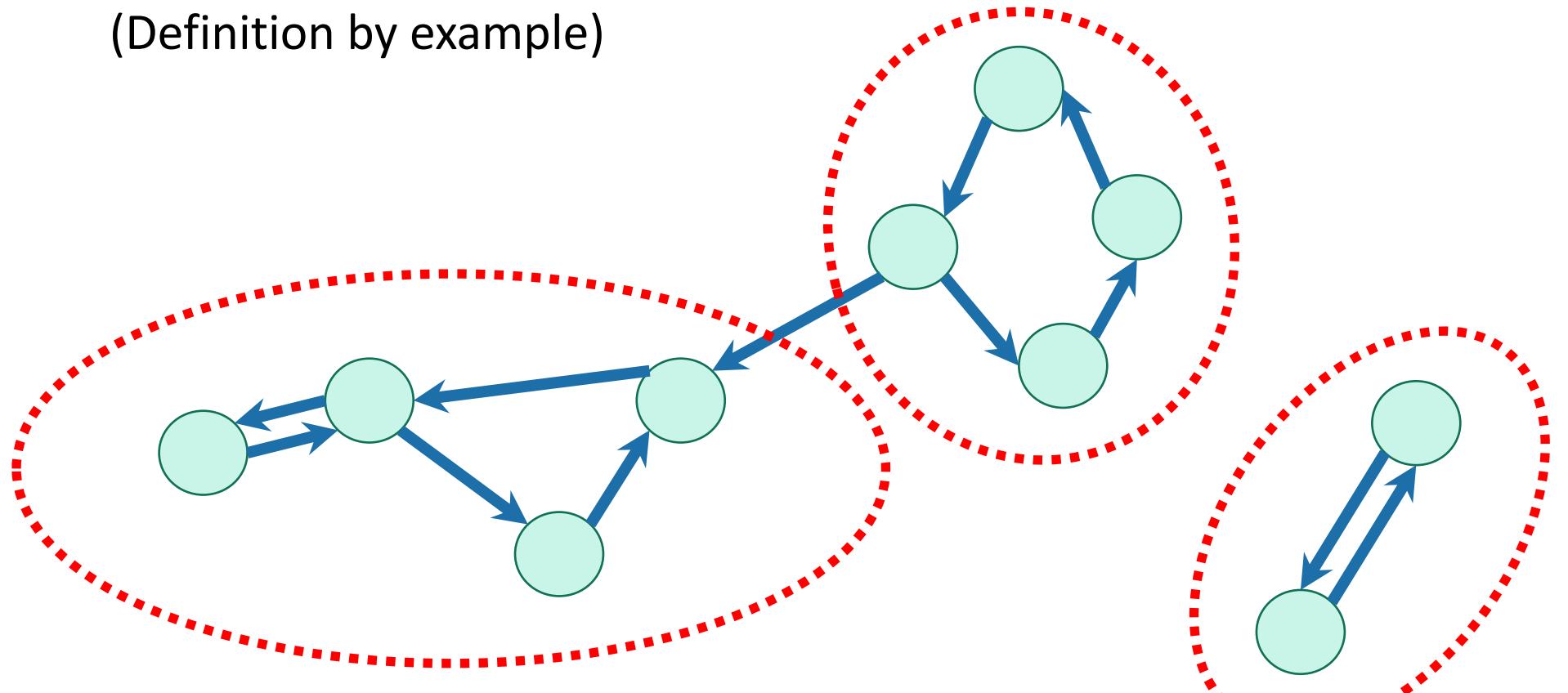
strongly connected



not strongly connected

We can decompose a graph into strongly connected components (SCCs)

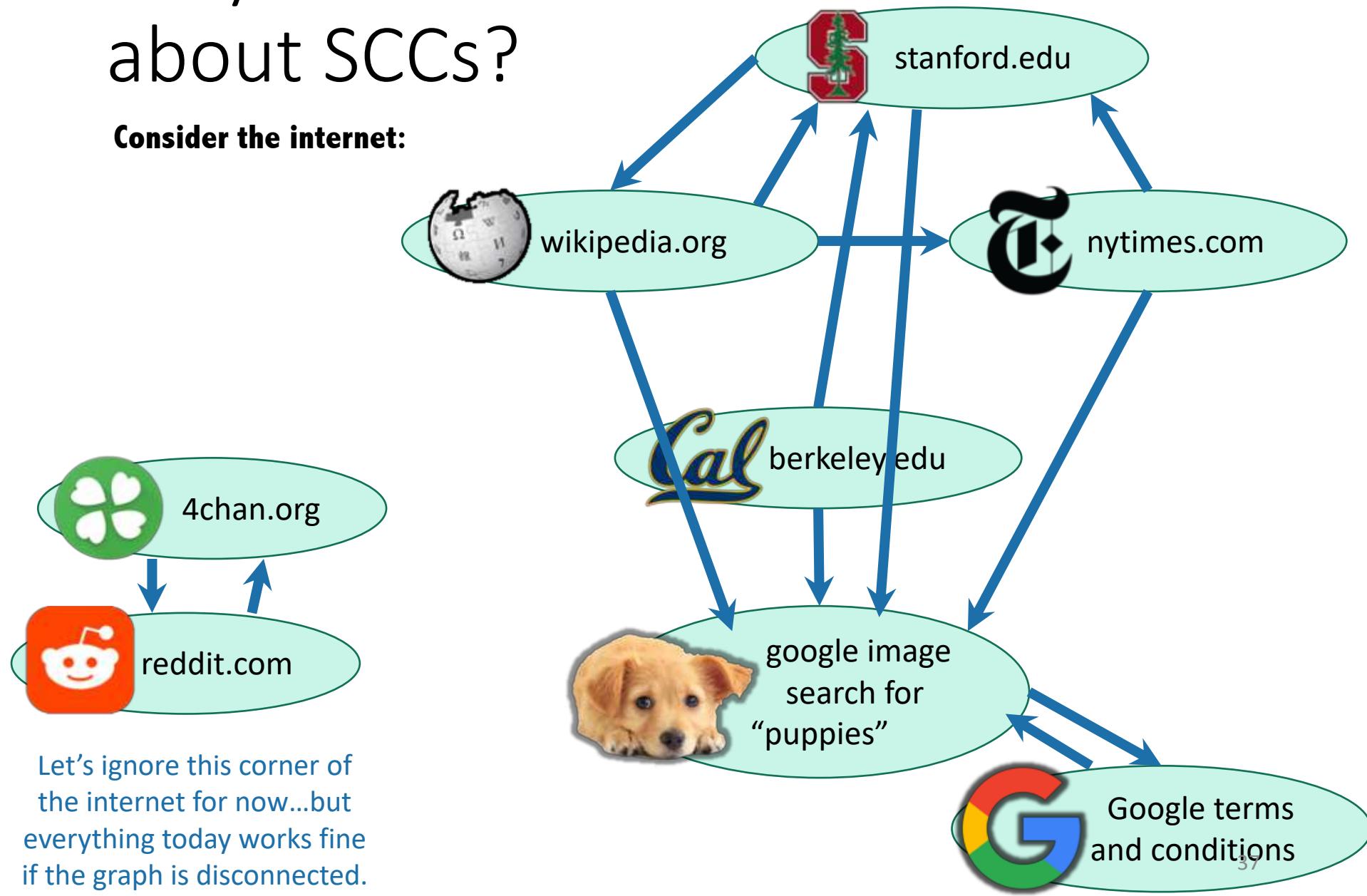
(Definition by example)



Definition by definition: The SCCs are the equivalence classes under the “are mutually reachable” equivalence relation.

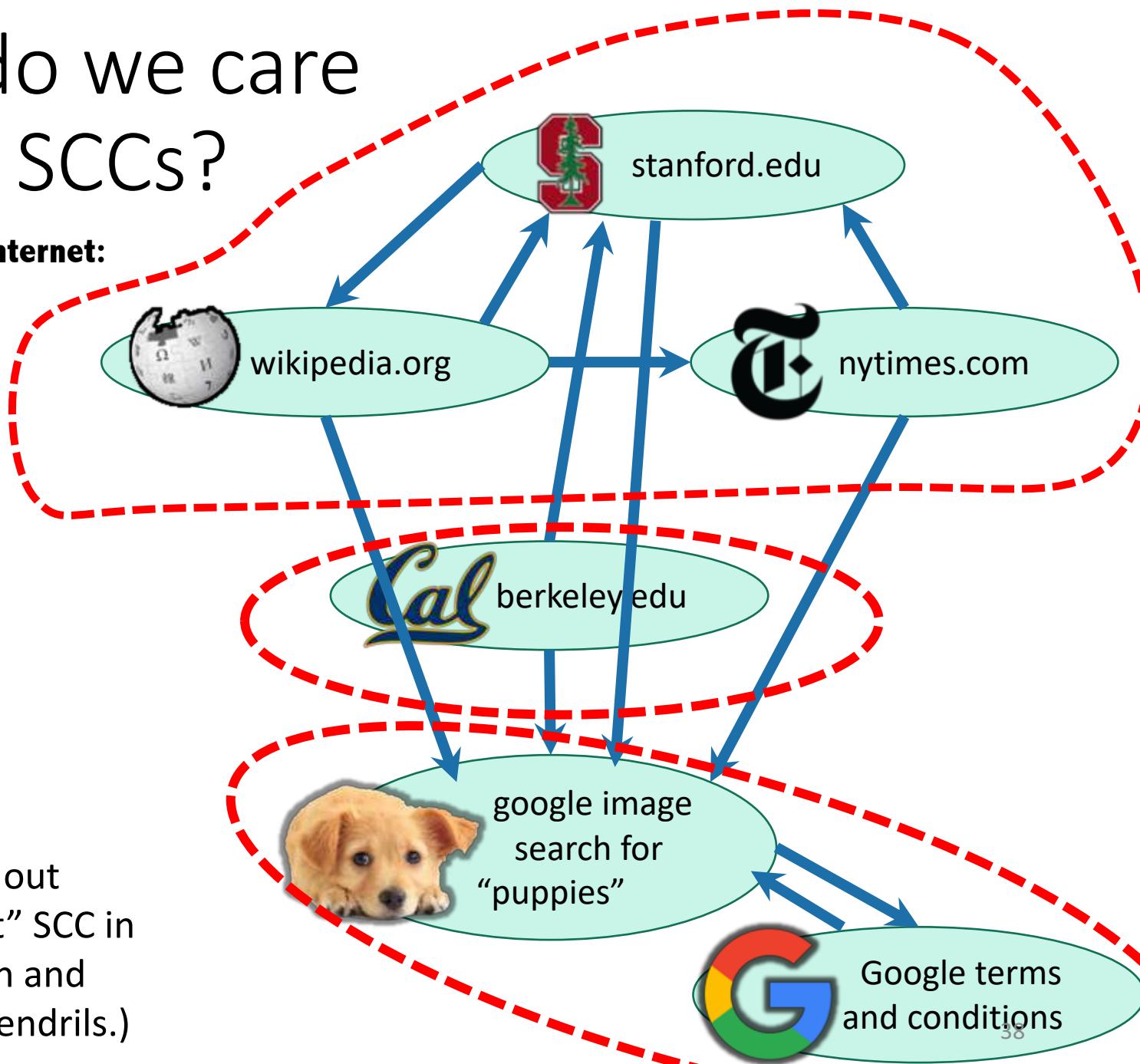
Why do we care about SCCs?

Consider the internet:



Why do we care about SCCs?

Consider the internet:



(In real life, turns out there's one “giant” SCC in the internet graph and then a bunch of tendrils.)

Why do we care about SCCs?

- Strongly connected components tell you about **communities**.
- Lots of graph algorithms only make sense on SCCs.
 - (So sometimes we want to find the SCCs as a first step)
 - Eg: I know an economist who has to first break up his labor market data into SCCs in order to make sense of it.

How to find SCCs?

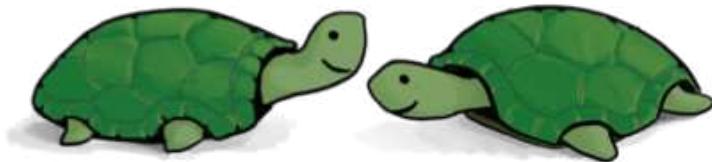
Try 1:

- Consider all possible decompositions and check.

Try 2:

- Something like...
 - Run DFS a bunch to find out which u's and v's belong in the same SCC
 - Aggregate that information to figure out the SCCs

Come up with a straightforward way to use DFS
to find SCCs. What's the running time?
More than n^2 or less than n^2 ?



One straightforward solution

- SCCs = []
- For each u:
 - Run DFS from u
 - For each vertex v that u can reach:
 - If v is in an SCC we've already found:
 - Run DFS from v to see if you can reach u
 - If so, add u to v's SCC
 - Break
 - If we didn't break, create a new SCC which just contains u.

This will not be our final solution so don't worry too much about it...



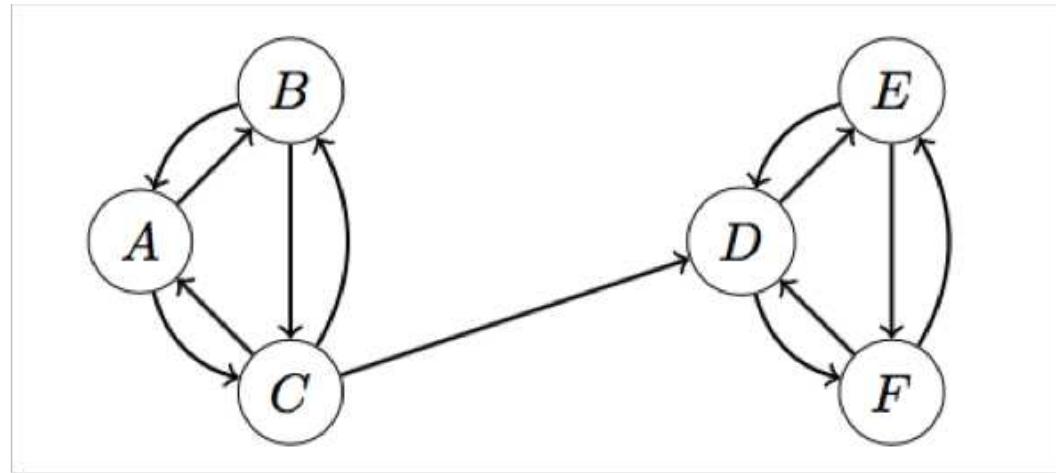
Running time AT LEAST $\Omega(n^2)$, no matter how smart you are about implementing the rest of it...

Today

- We will see how to find strongly connected components in time $O(n+m)$
- !!!!!

Pre-Lecture exercise

- Run DFS starting at D:



- That will identify SCCs...
- Issues:
 - How do we know where to start DFS?
 - It wouldn't have found the SCCs if we started from A.

Algorithm

Running time: $O(n + m)$

- Do DFS to create a DFS forest.
 - Choose starting vertices in any order.
 - Keep track of finishing times.
- Reverse all the edges in the graph.
- Do DFS again to create **another DFS forest**.
 - This time, order the nodes in the reverse order of the finishing times that they had from the first DFS run.
- The SCCs are the different trees in the **second DFS forest**.



Look, it works!

- (See IPython notebook)

```
In [4]: print(G)
```

```
CS161Graph with:  
    Vertices:  
        Stanford, Wikipedia, NYTimes, Berkeley, Puppies, Google,  
    Edges:  
        (Stanford, Wikipedia) (Stanford, Puppies) (Wikipedia, Stanford) (Wikipedia, NYTimes)  
        (Wikipedia, Puppies) (NYTimes, Stanford) (NYTimes, Puppies)  
        (Berkeley, Stanford) (Berkeley, Puppies) (Puppies, Google) (Google, Puppies)
```

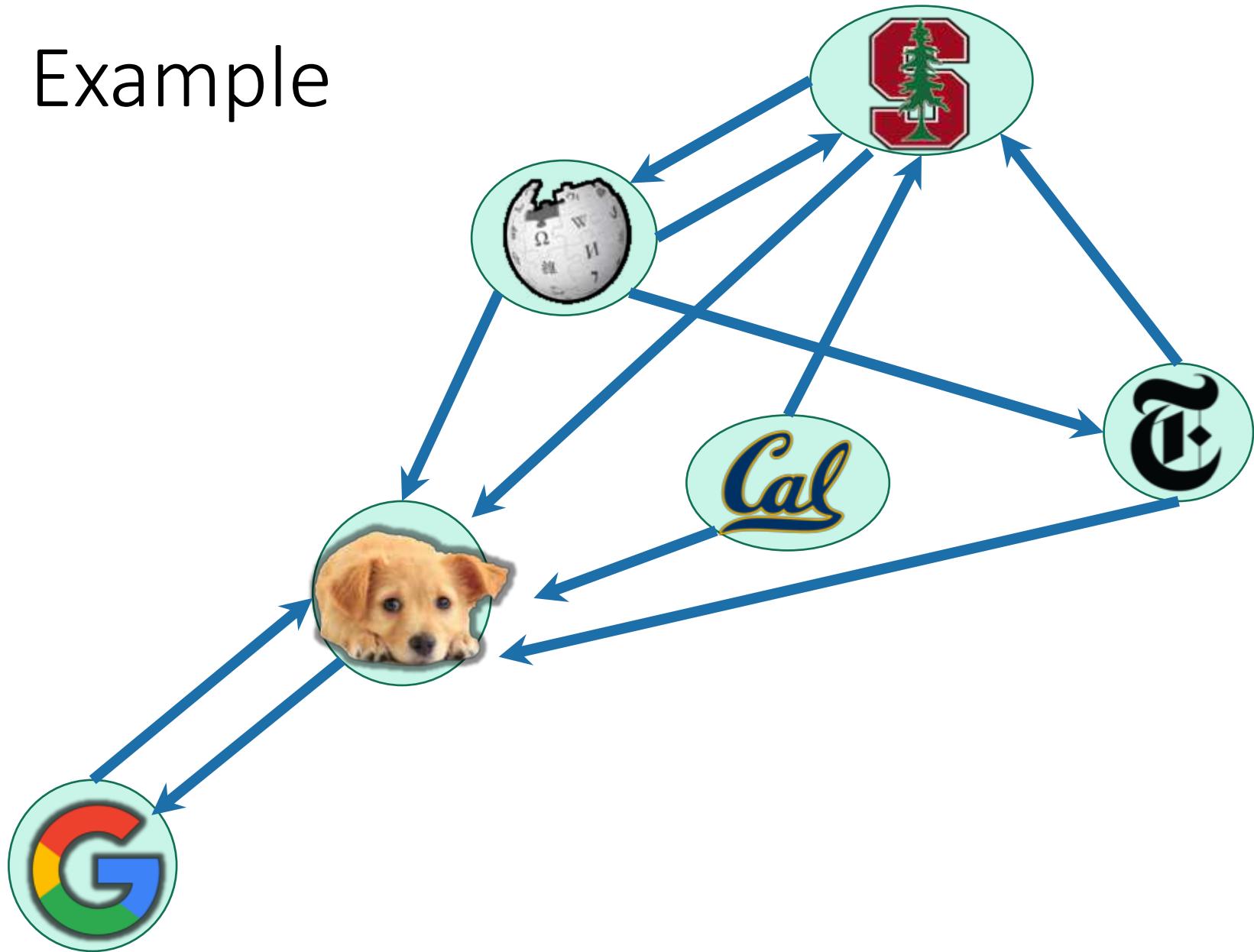
```
In [5]: SCCs = SCC(G, False)
```

```
for X in SCCs:  
    print ([str(x) for x in X])
```

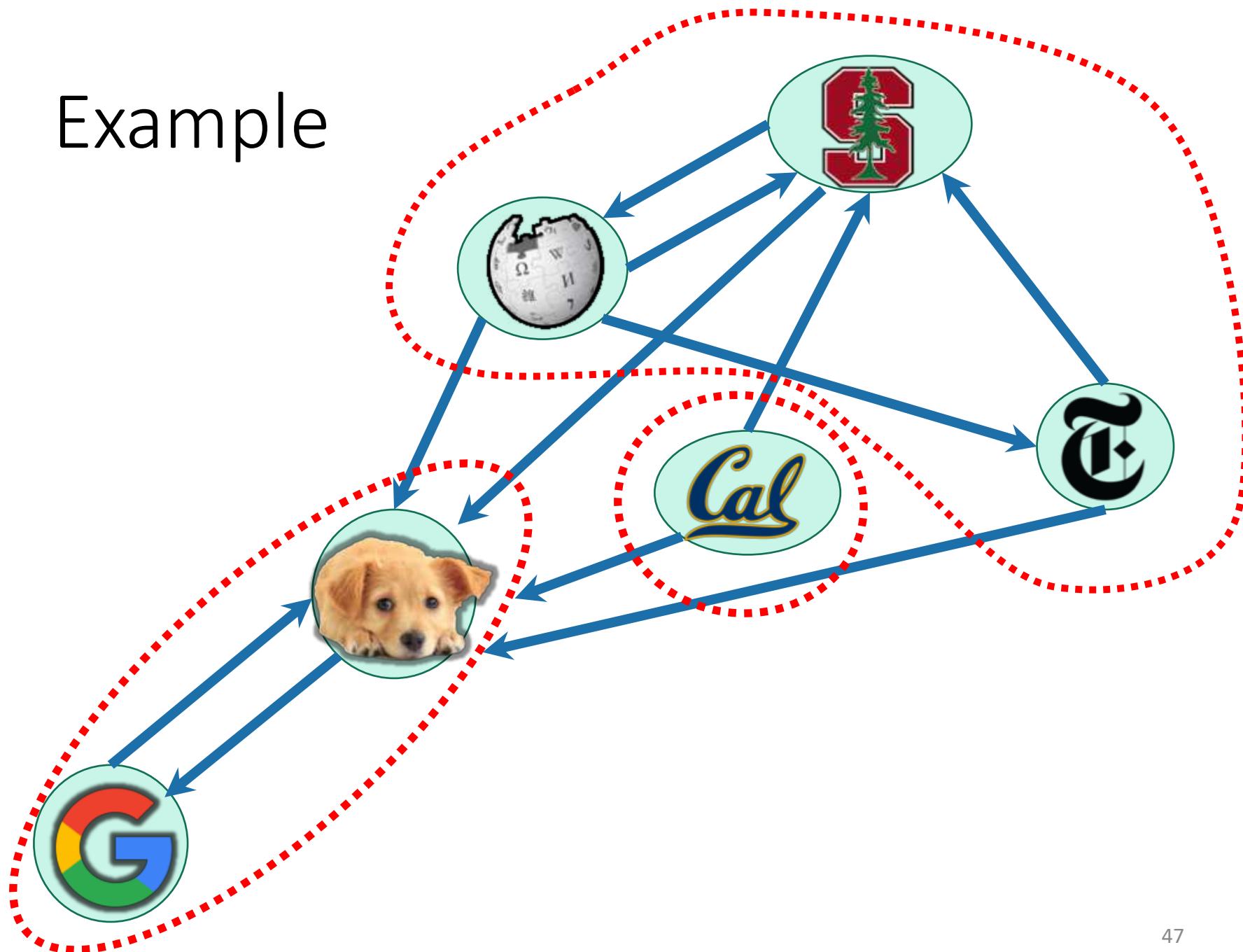
```
['Berkeley']  
['Stanford', 'NYTimes', 'Wikipedia']  
['Puppies', 'Google']
```

But let's break that down a bit...

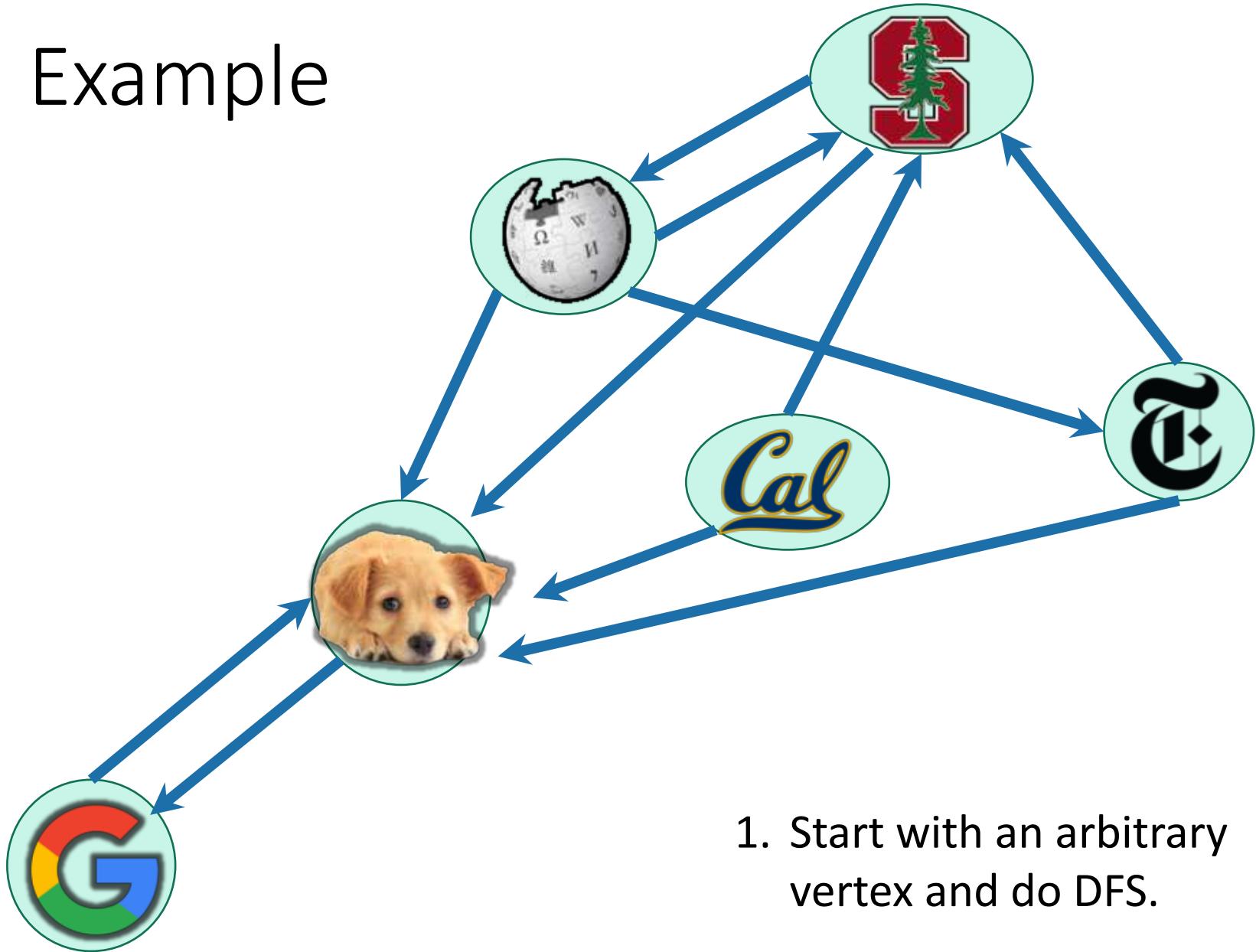
Example



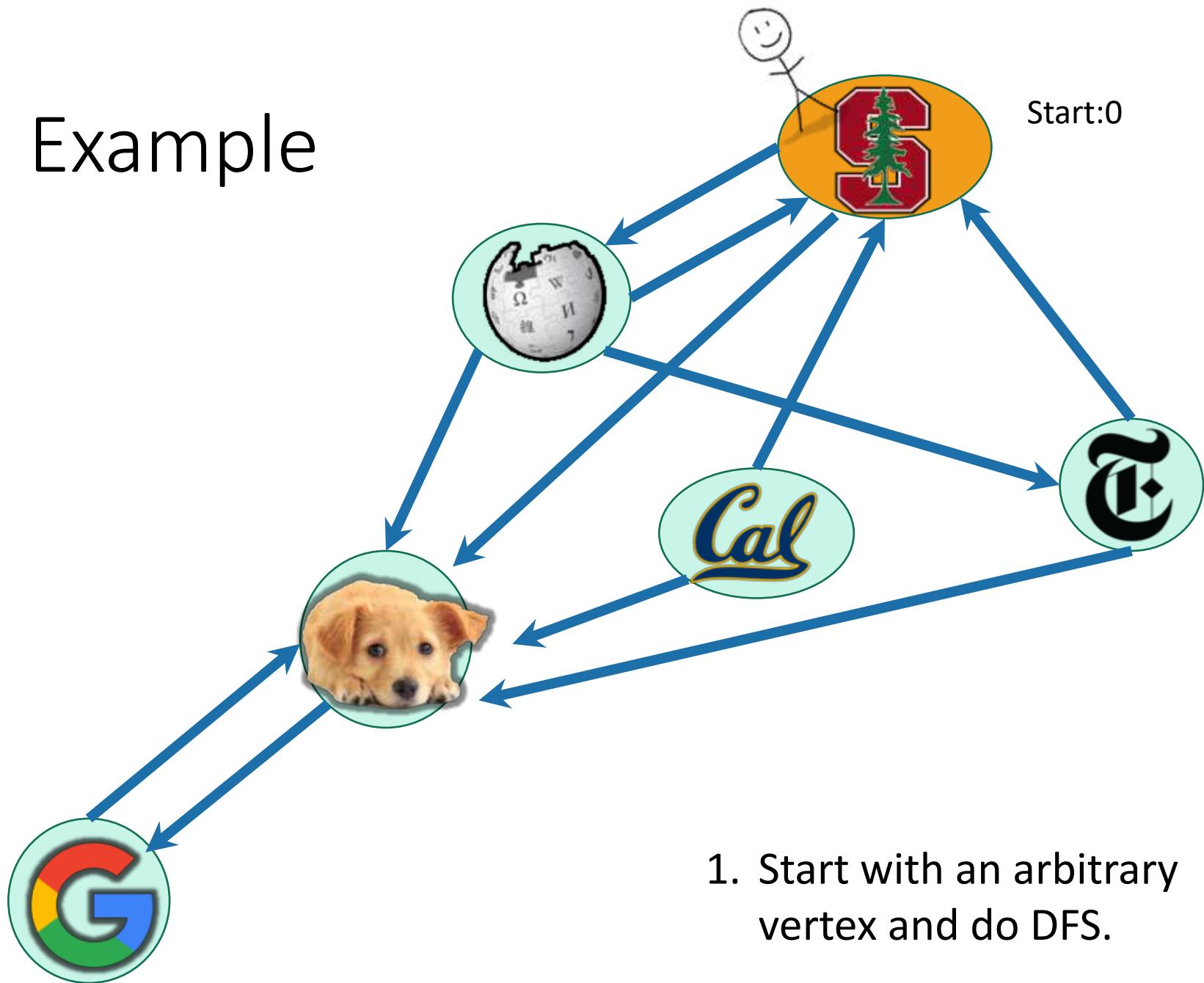
Example



Example

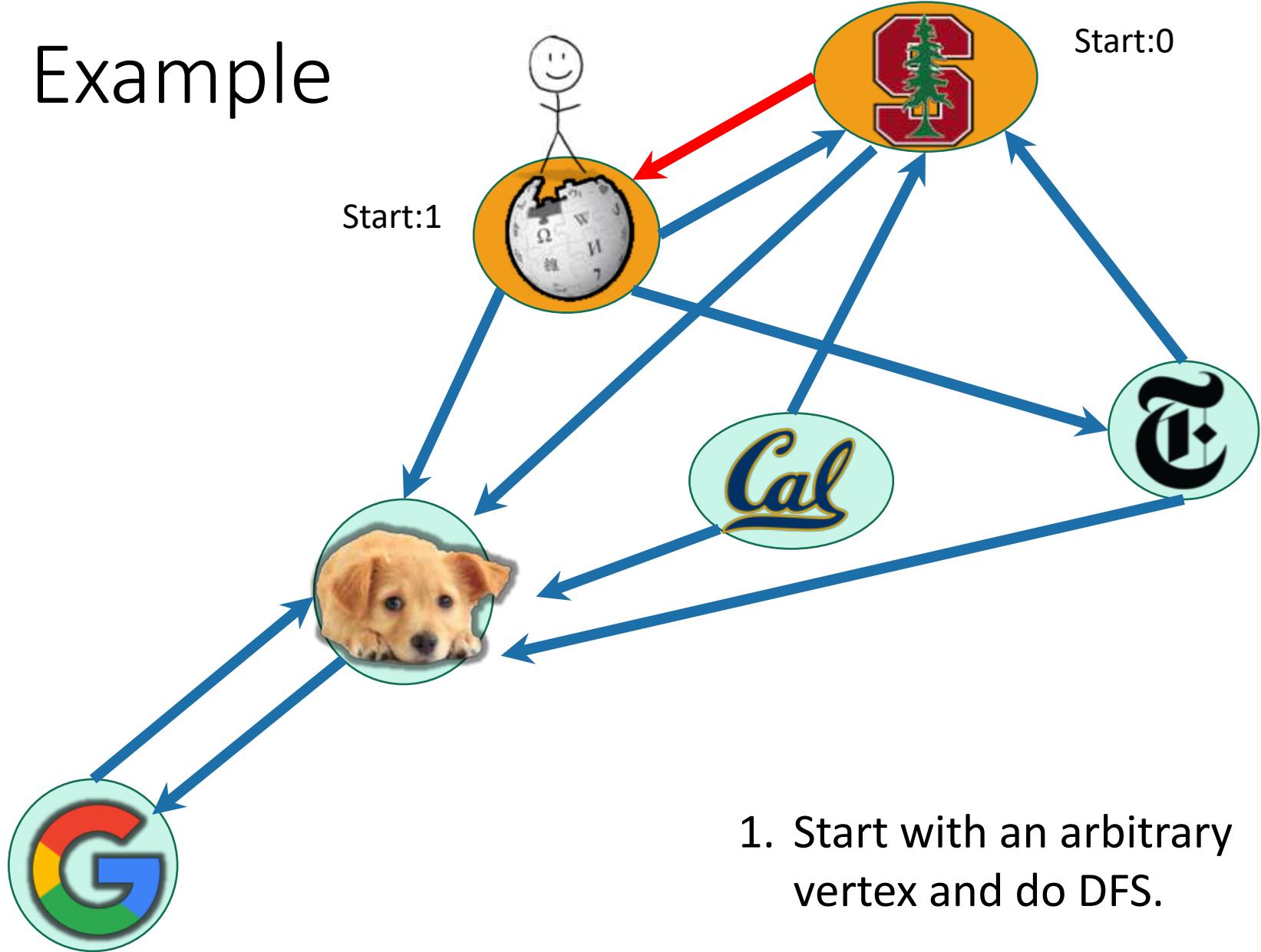


Example

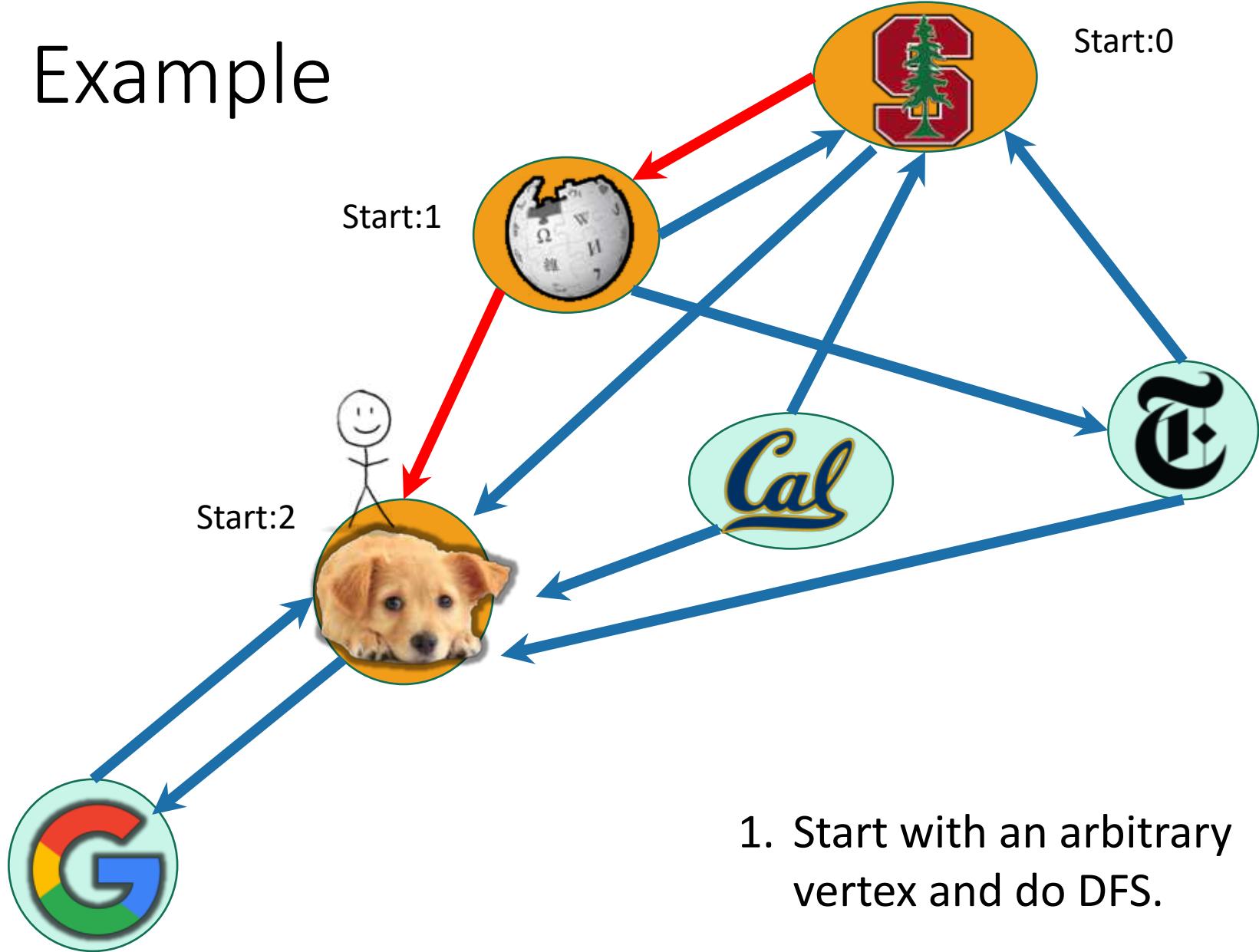


1. Start with an arbitrary vertex and do DFS.

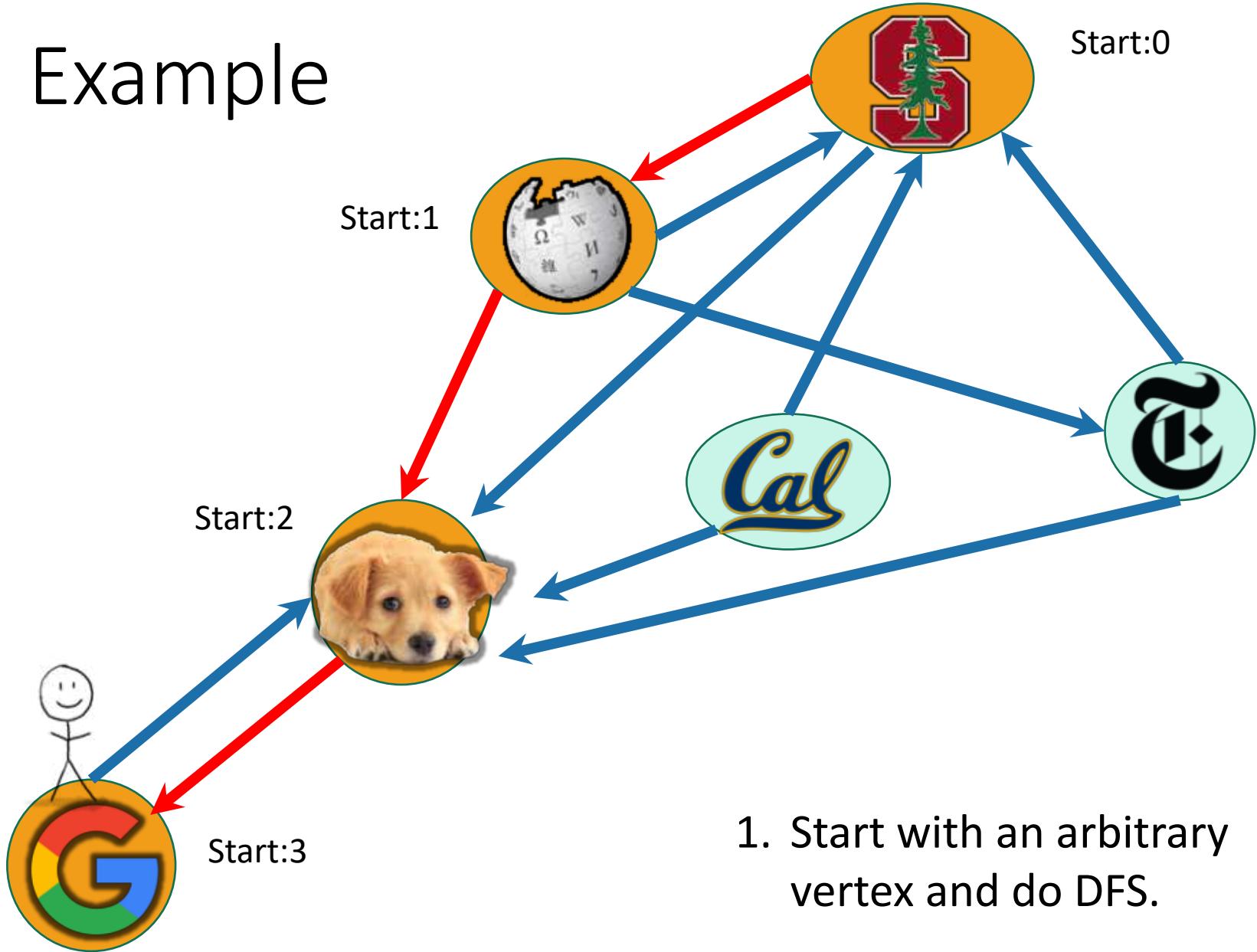
Example



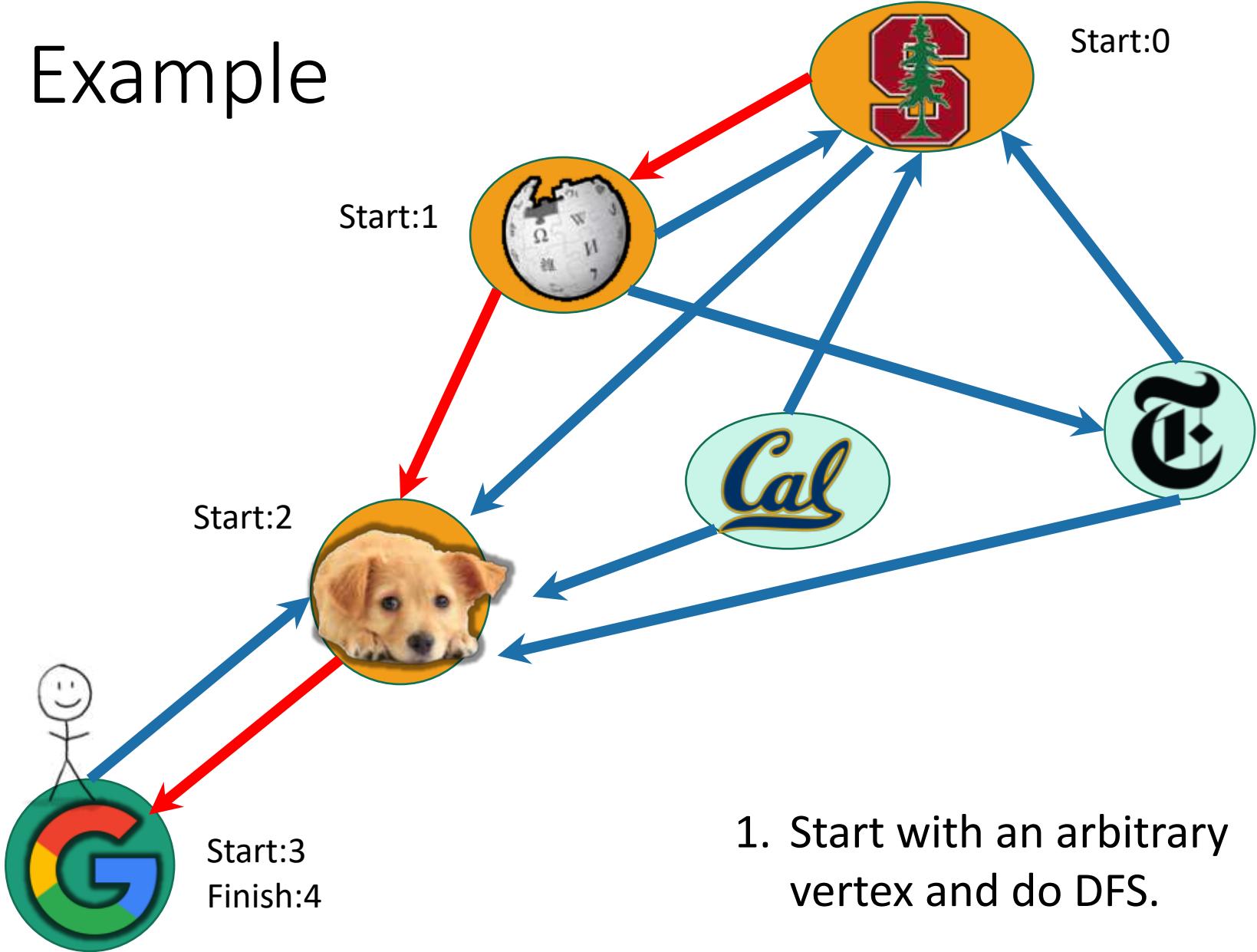
Example



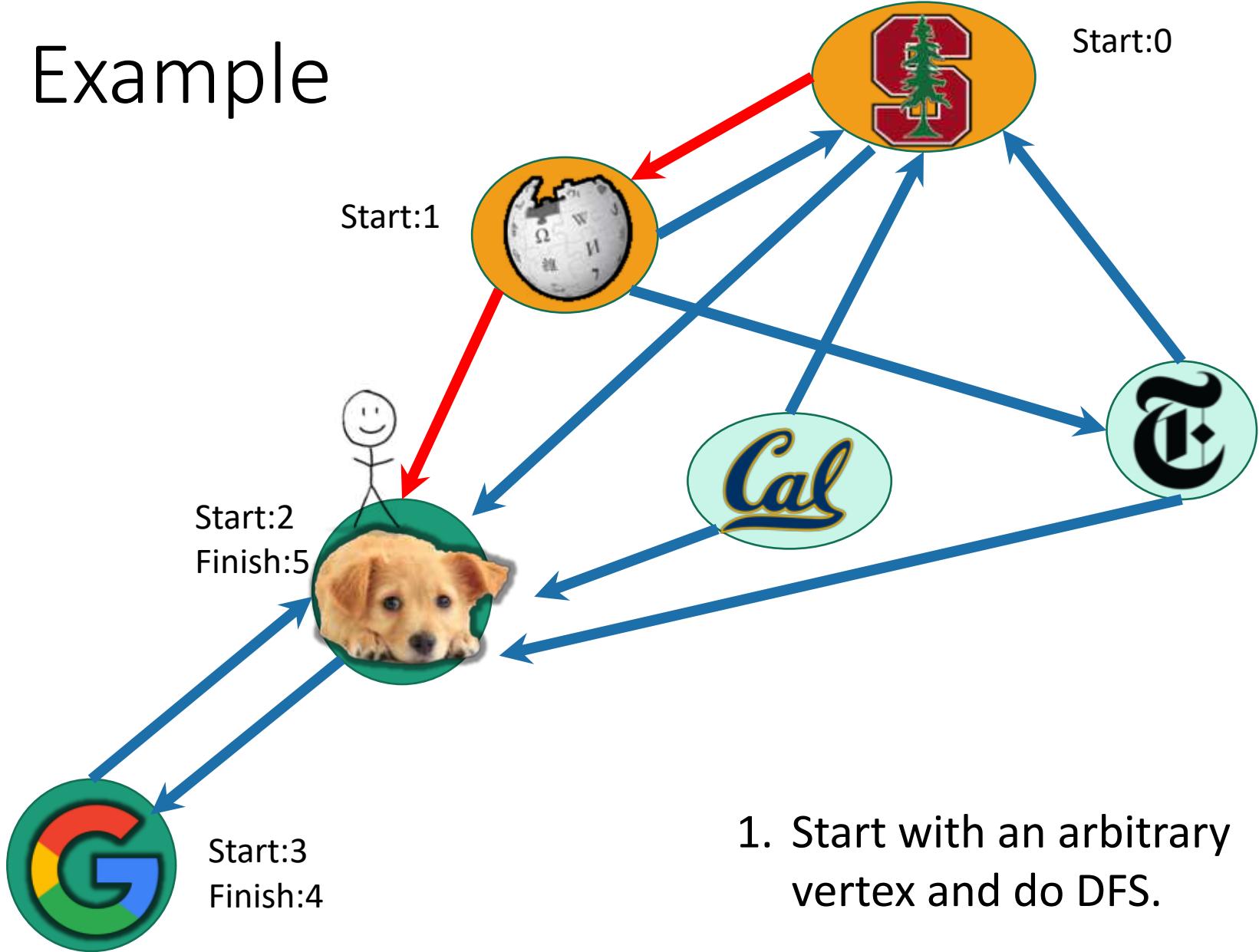
Example



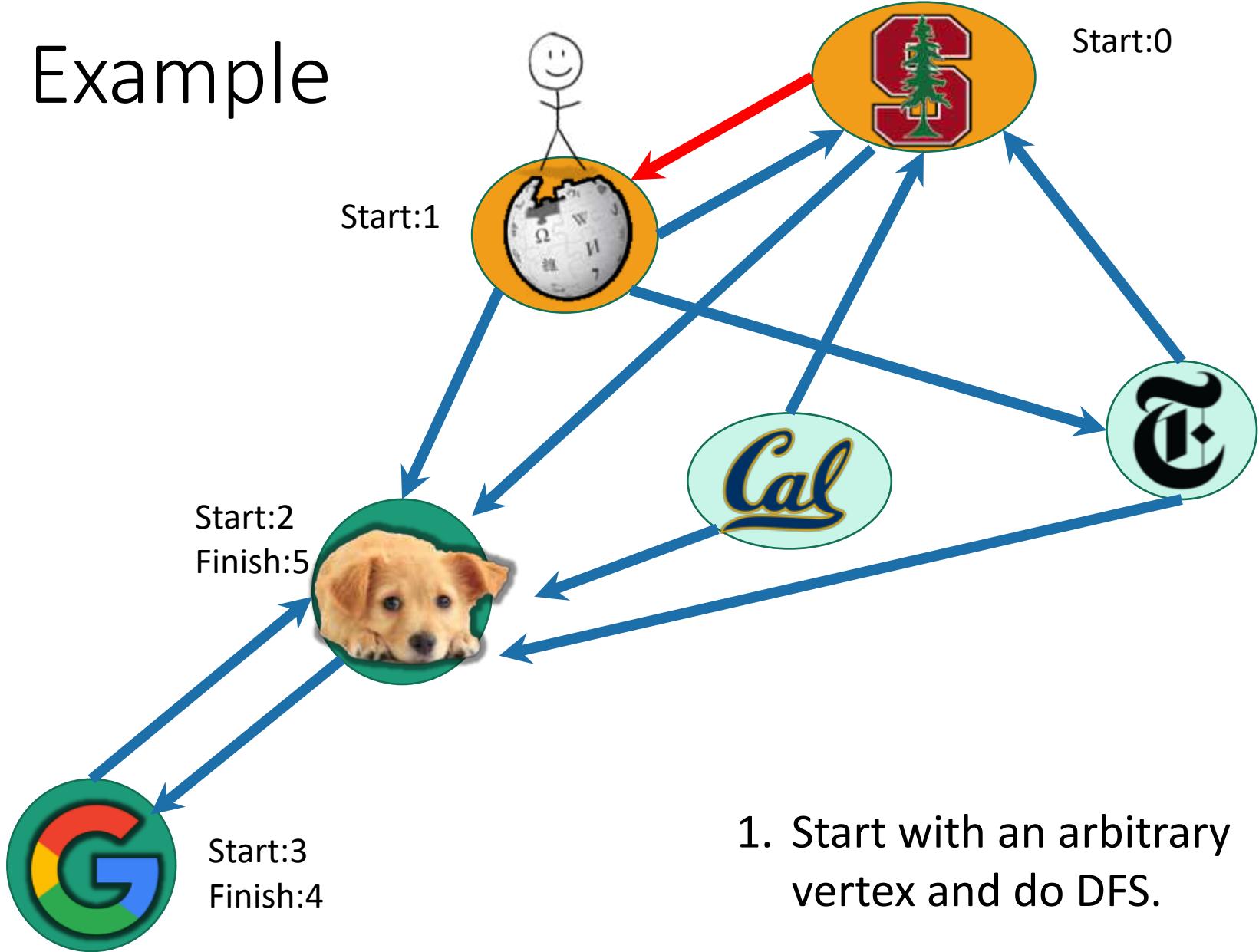
Example



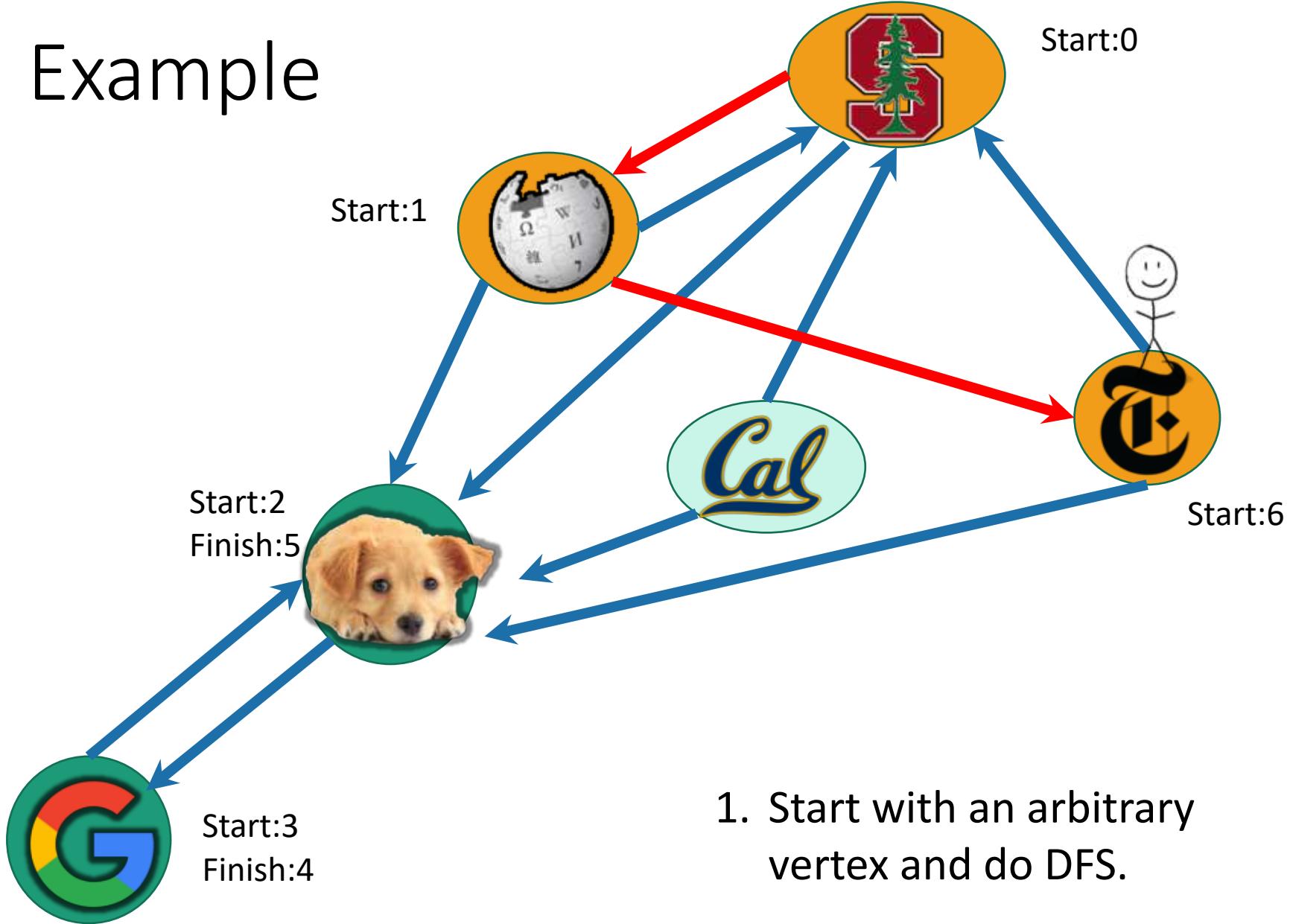
Example



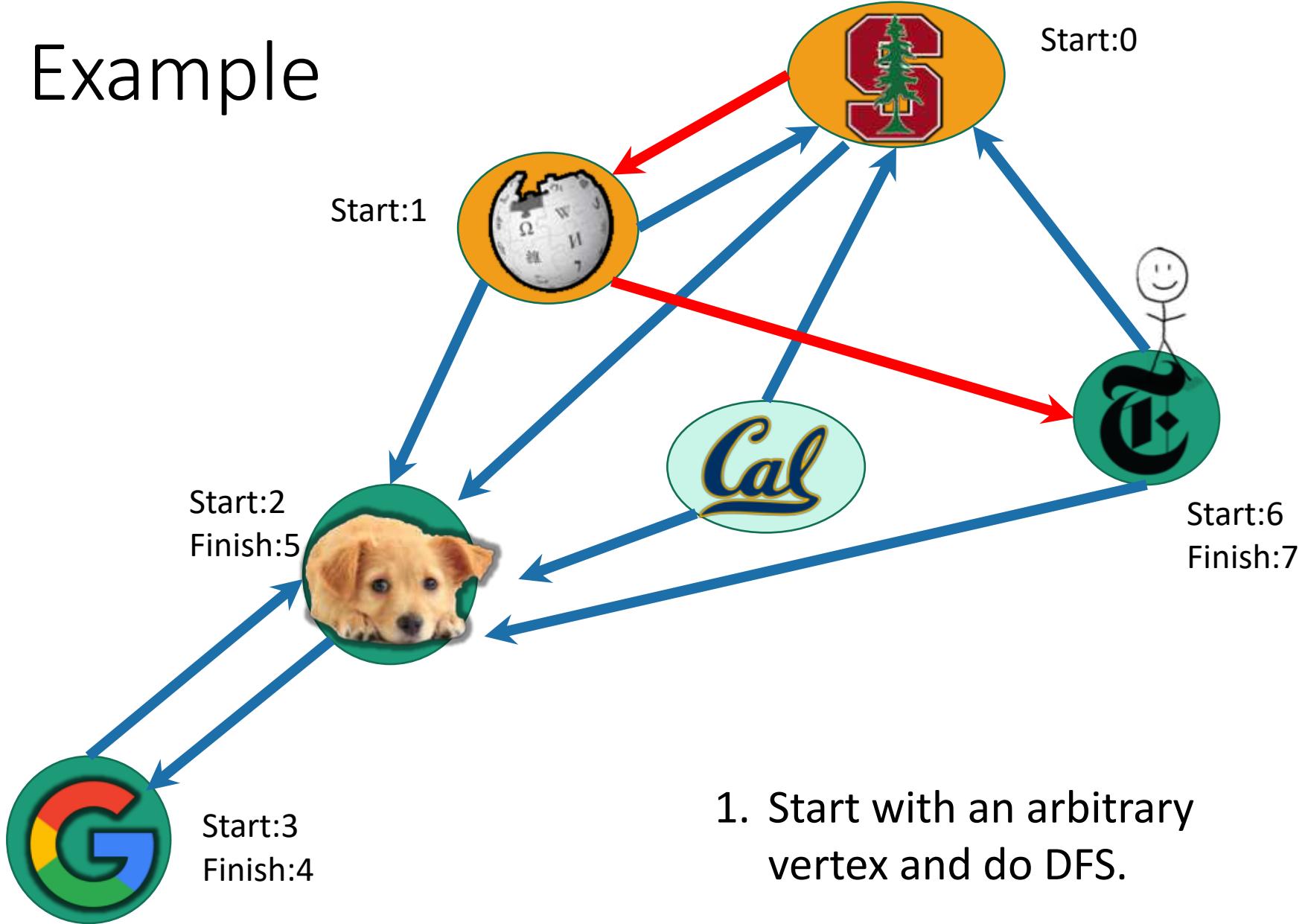
Example



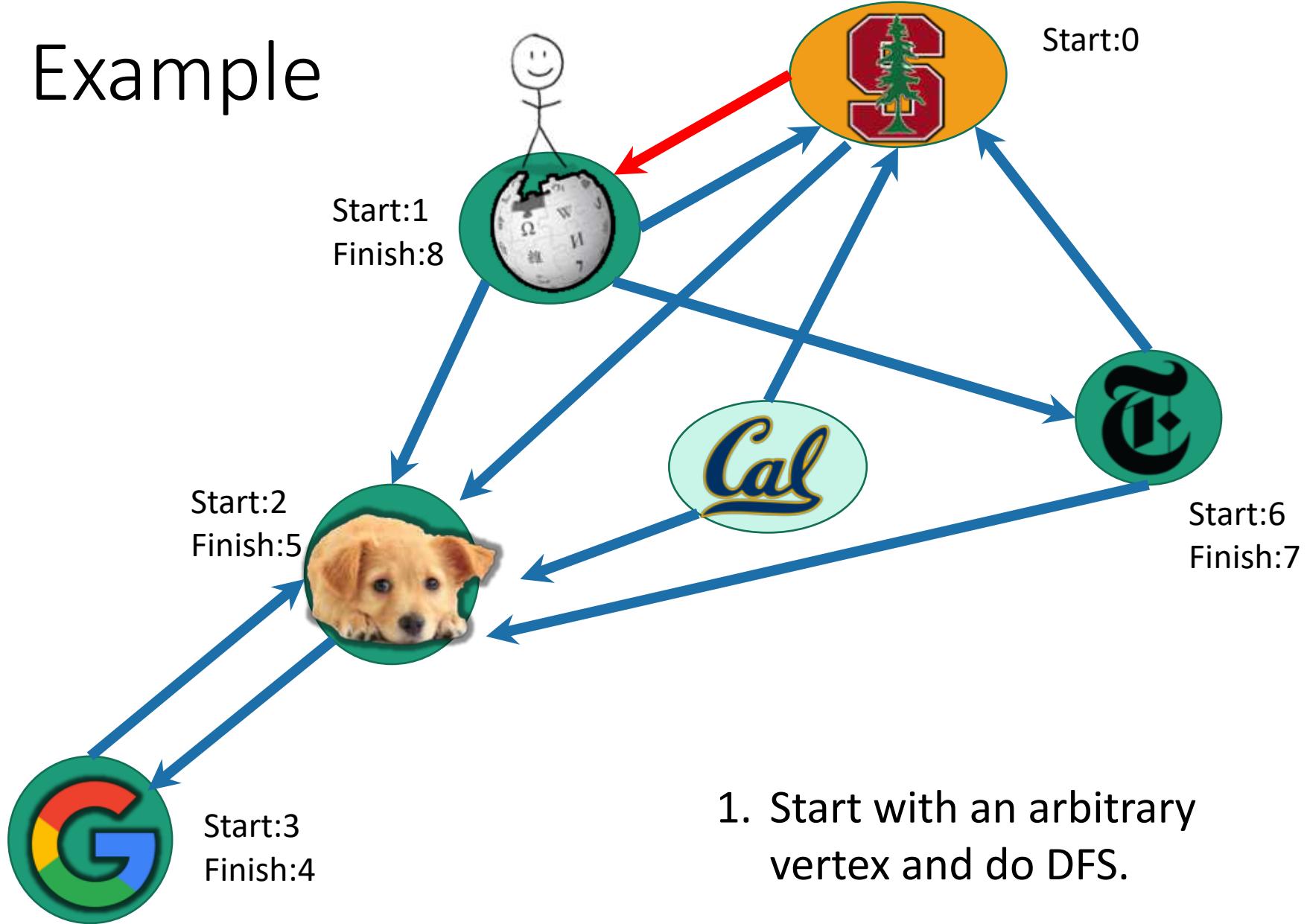
Example



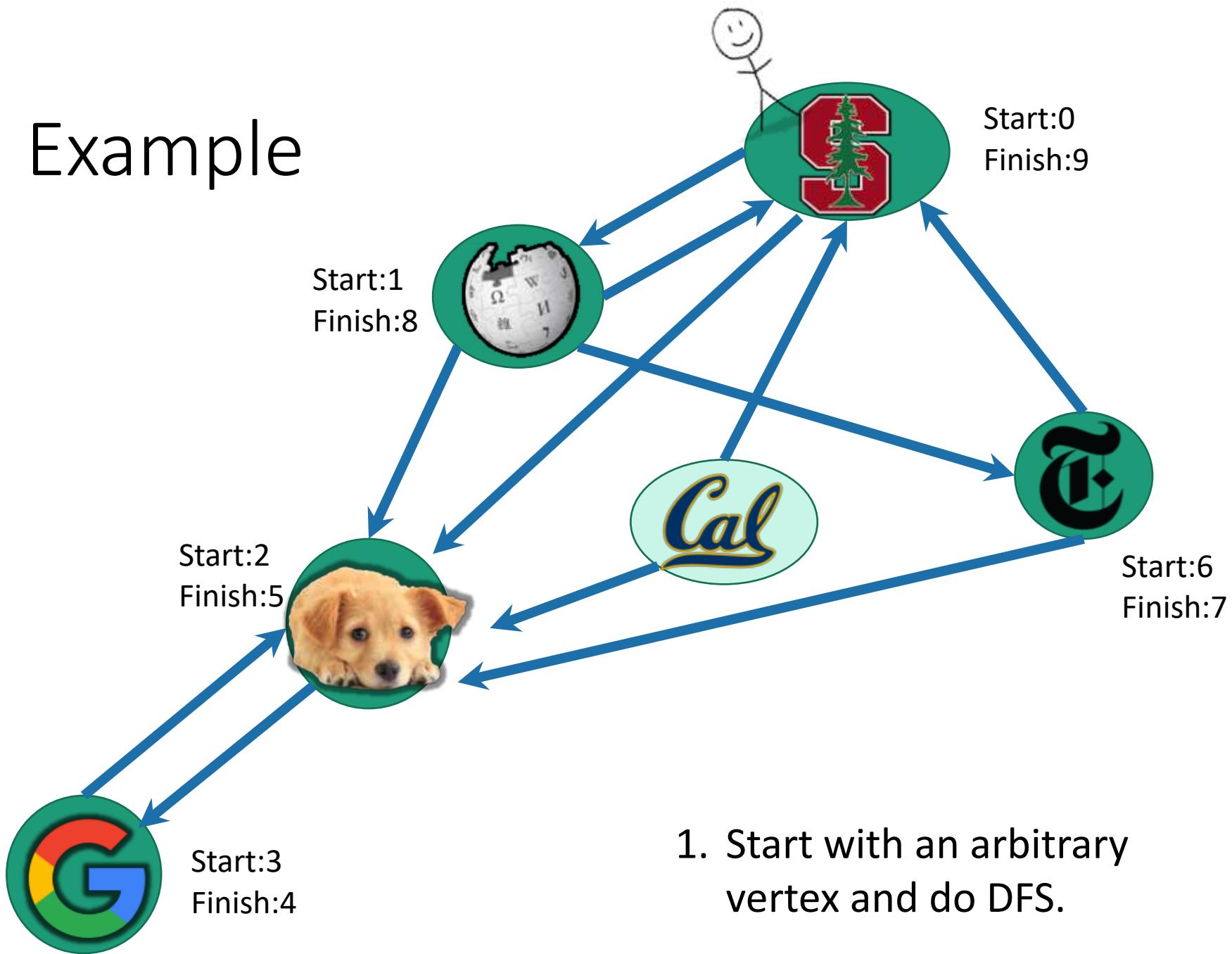
Example



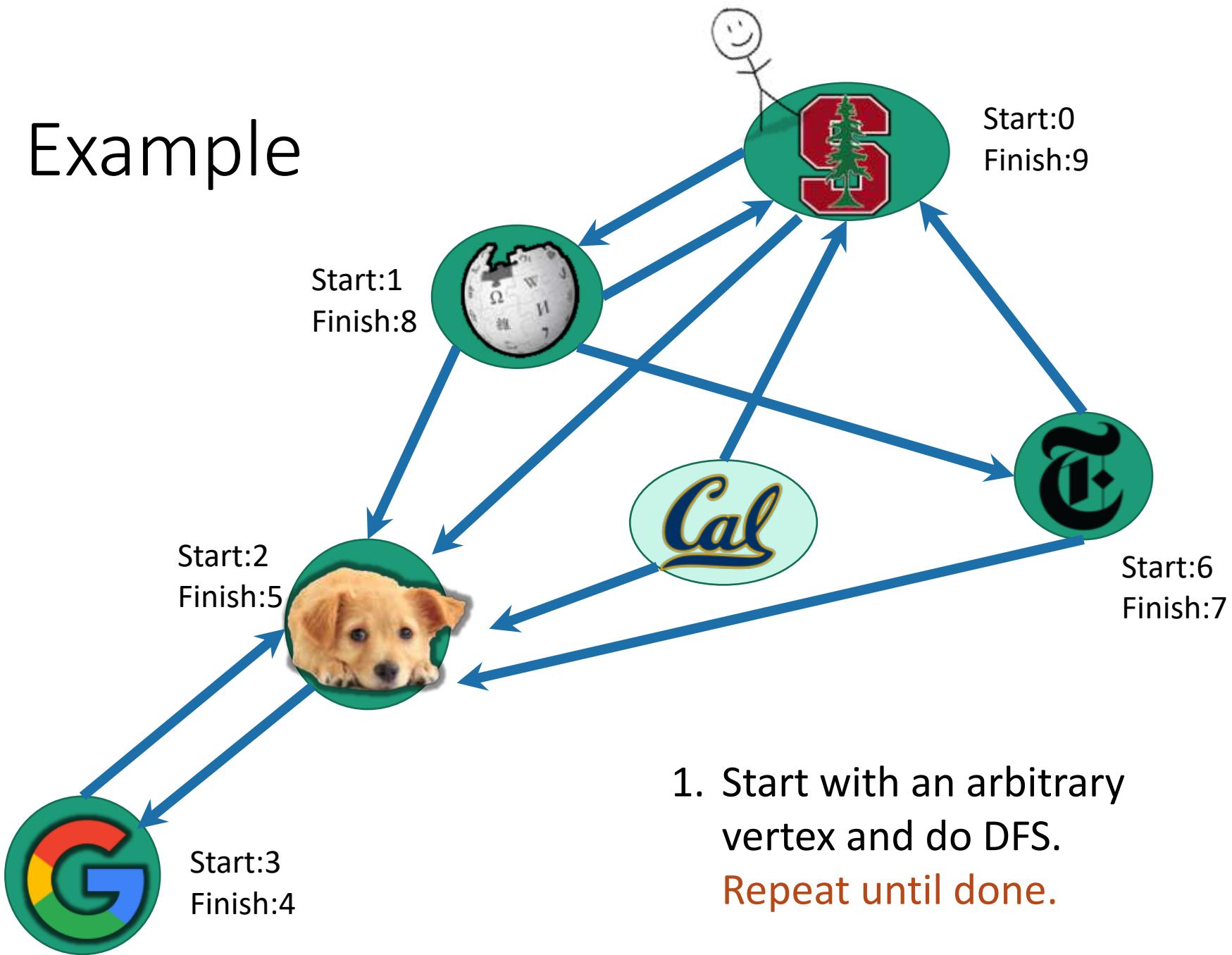
Example



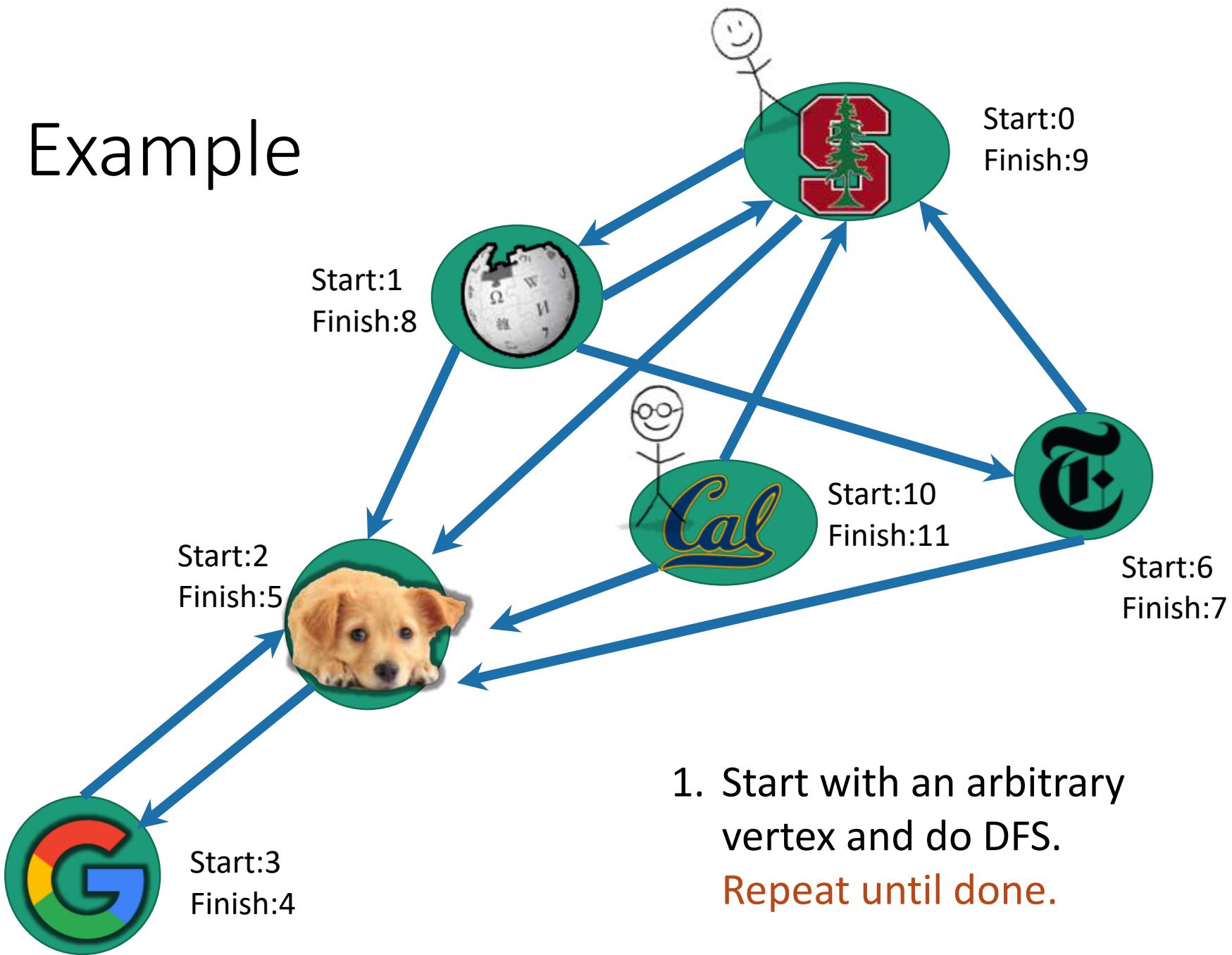
Example



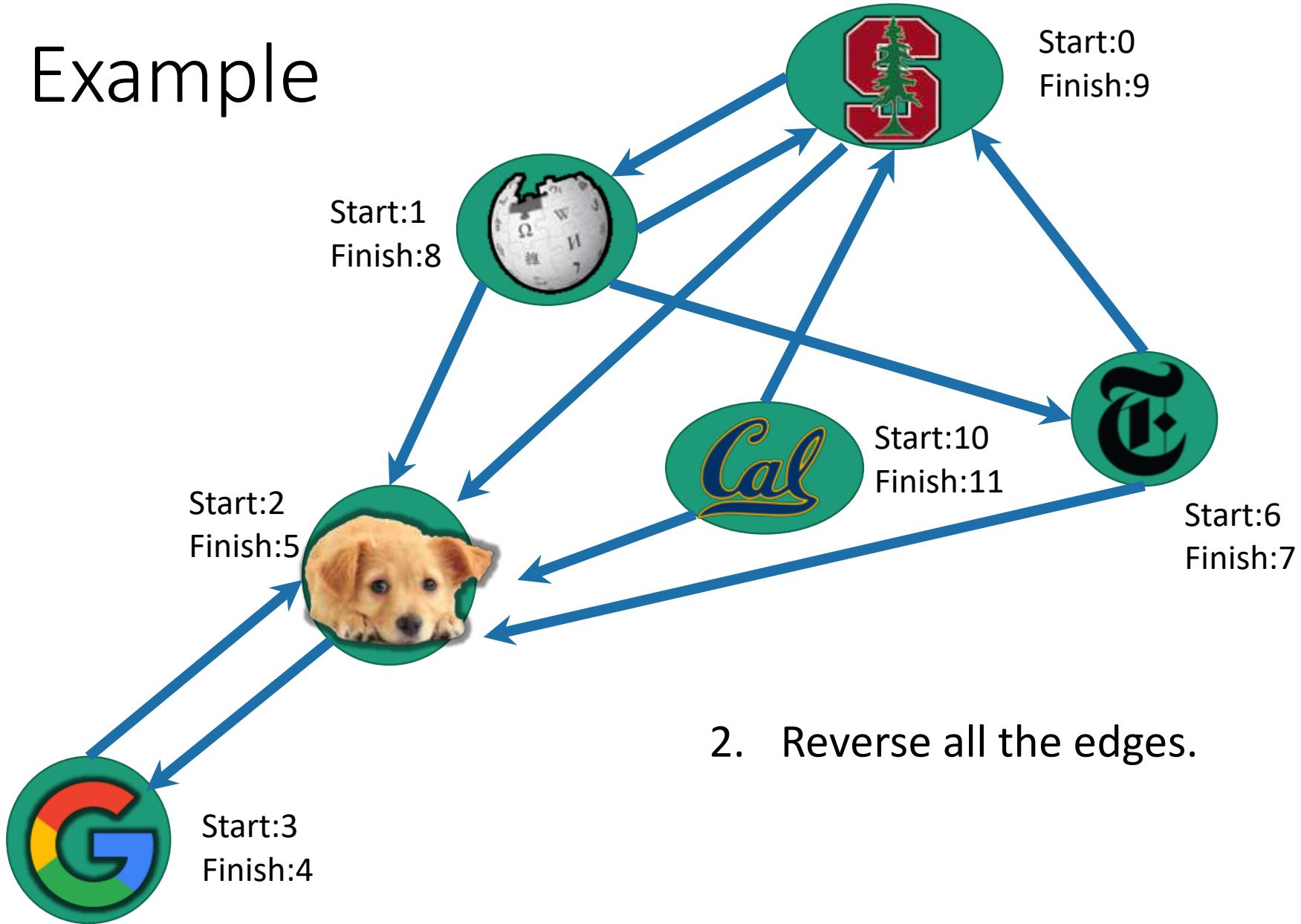
Example



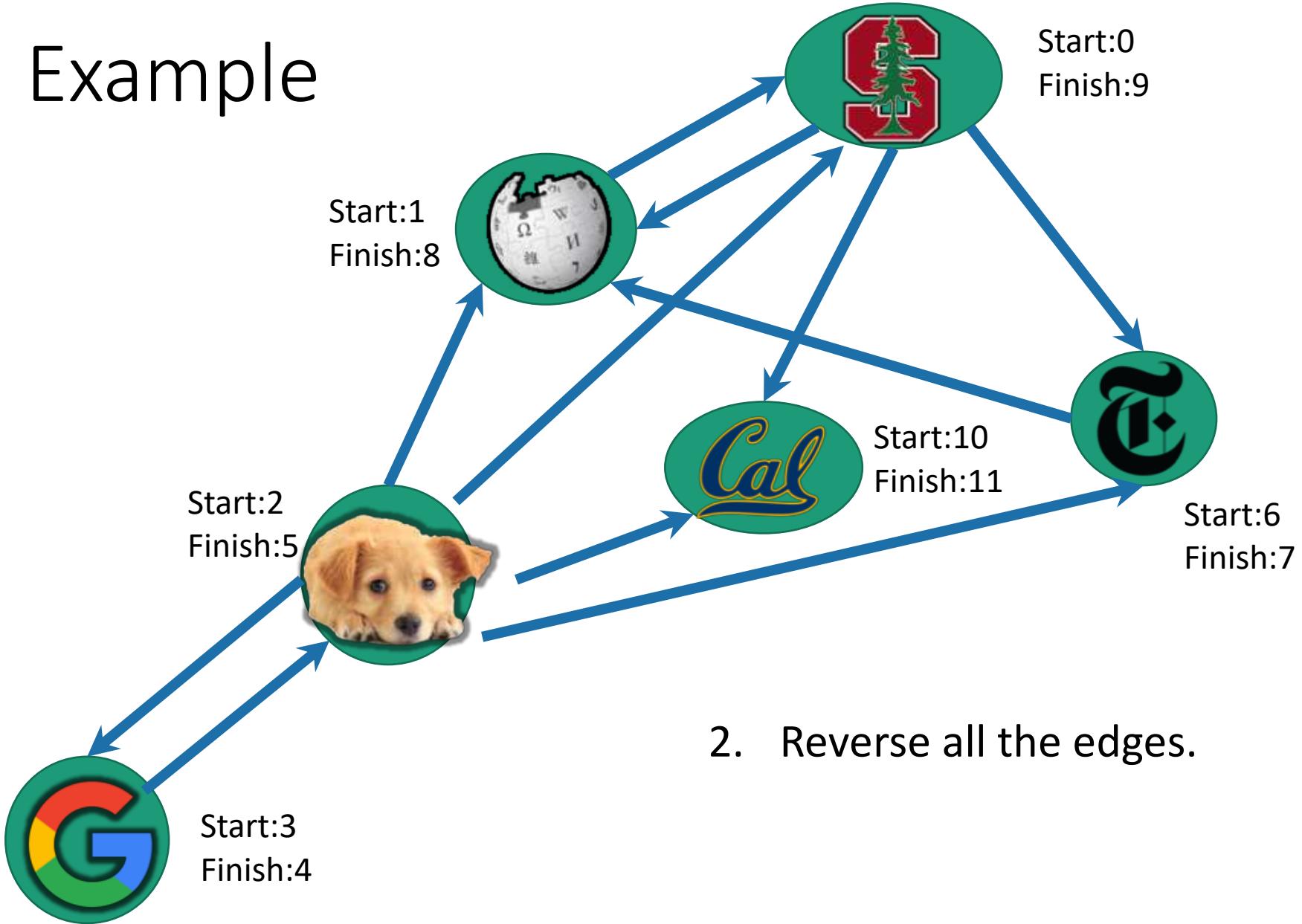
Example



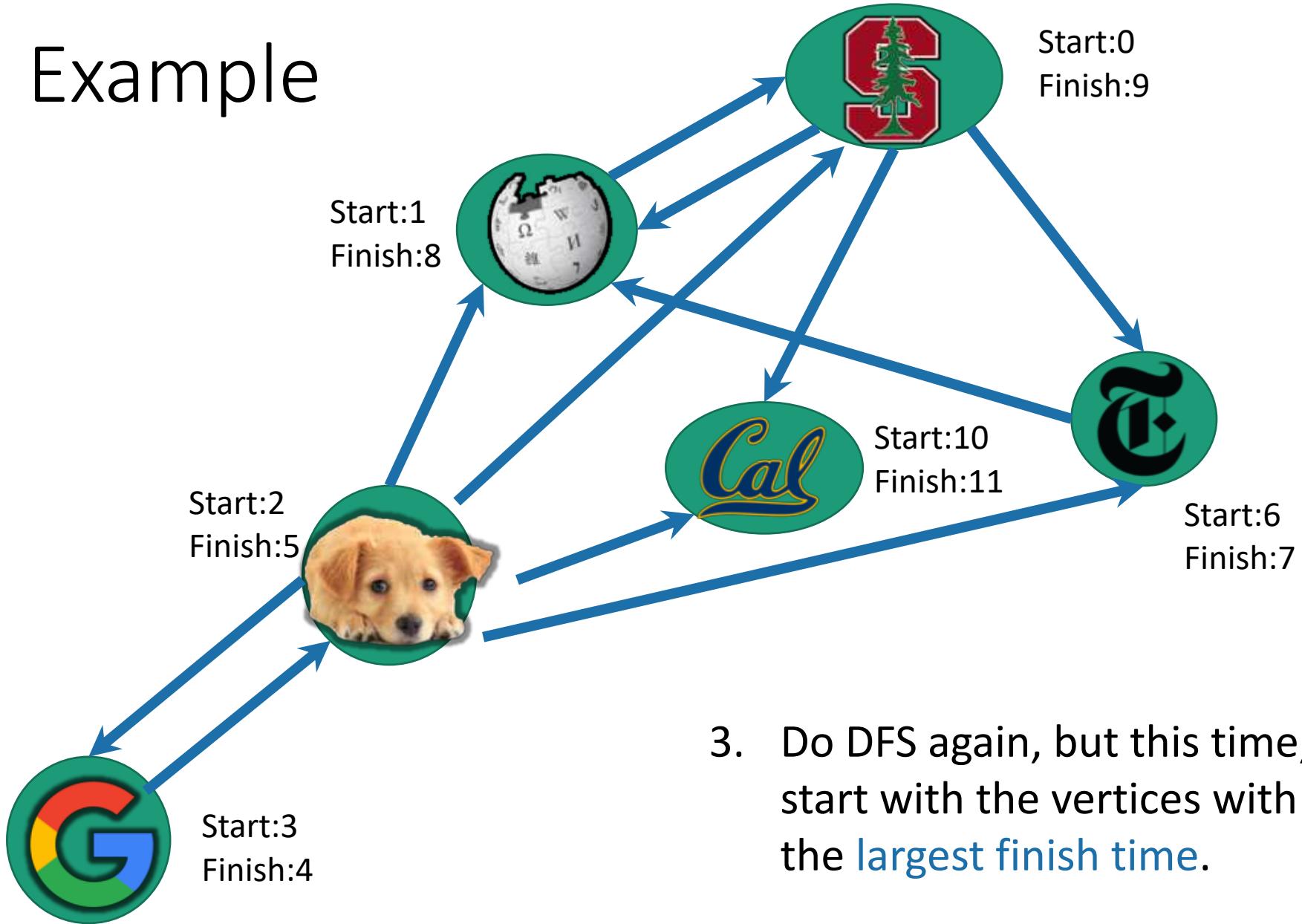
Example



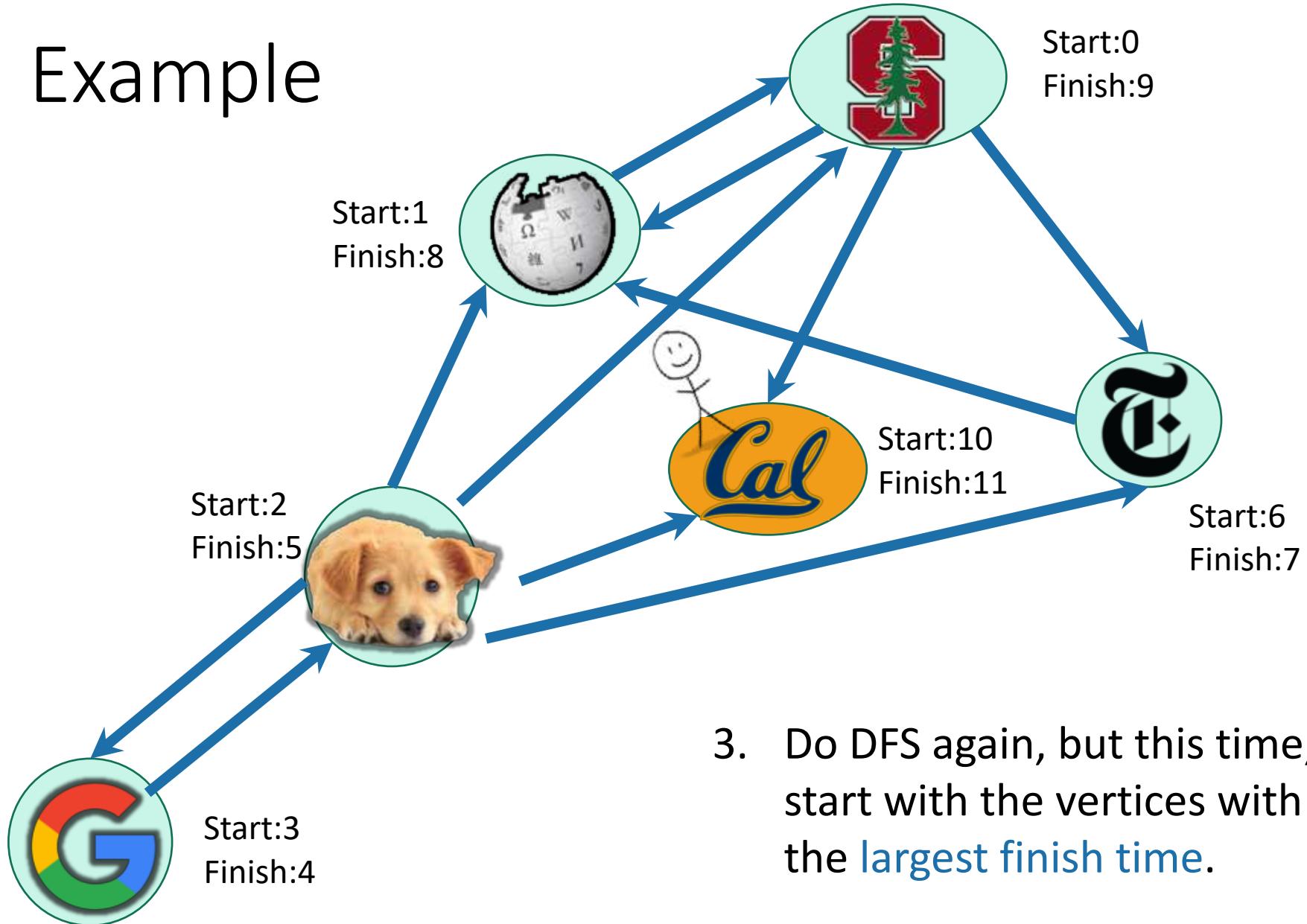
Example



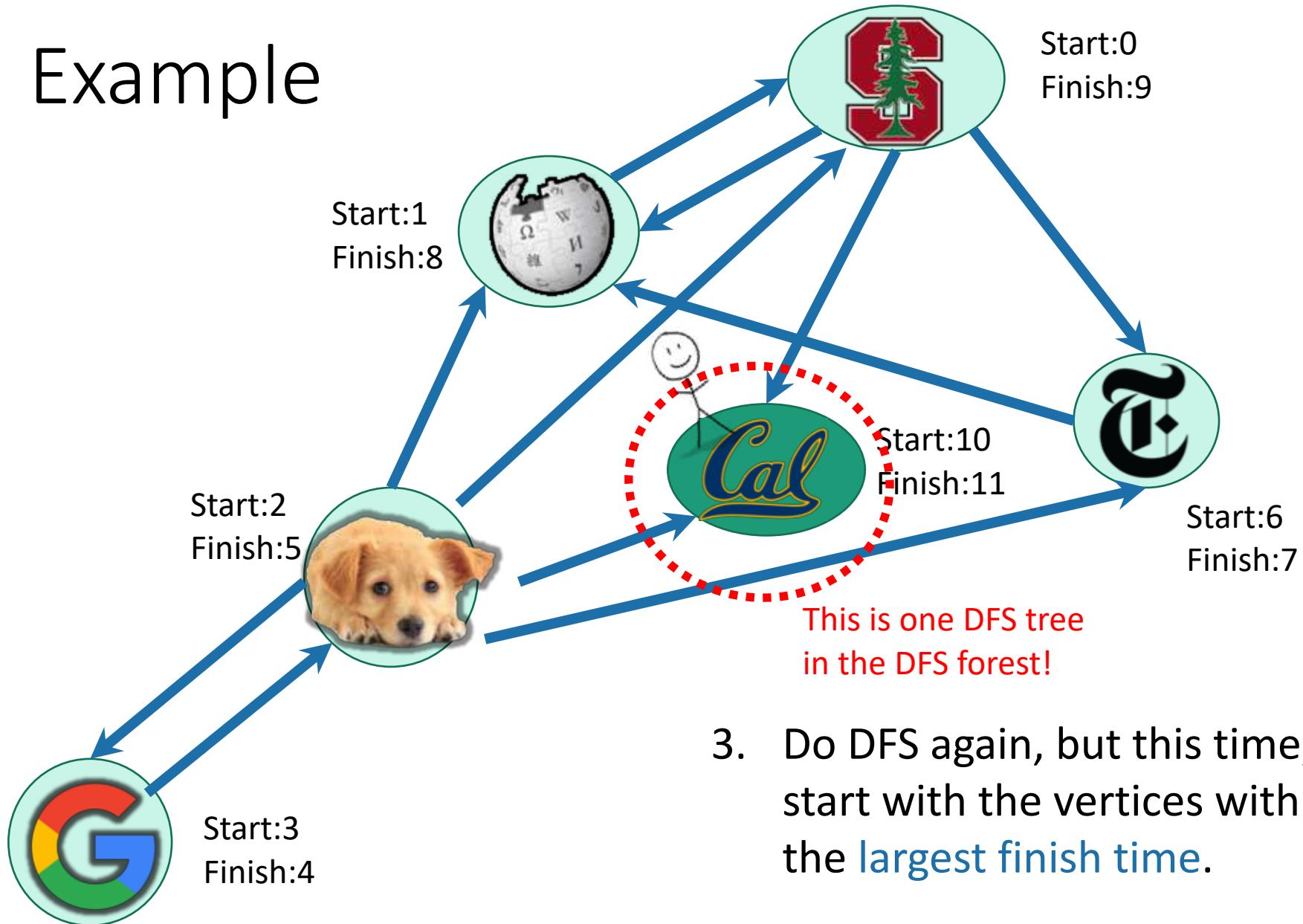
Example



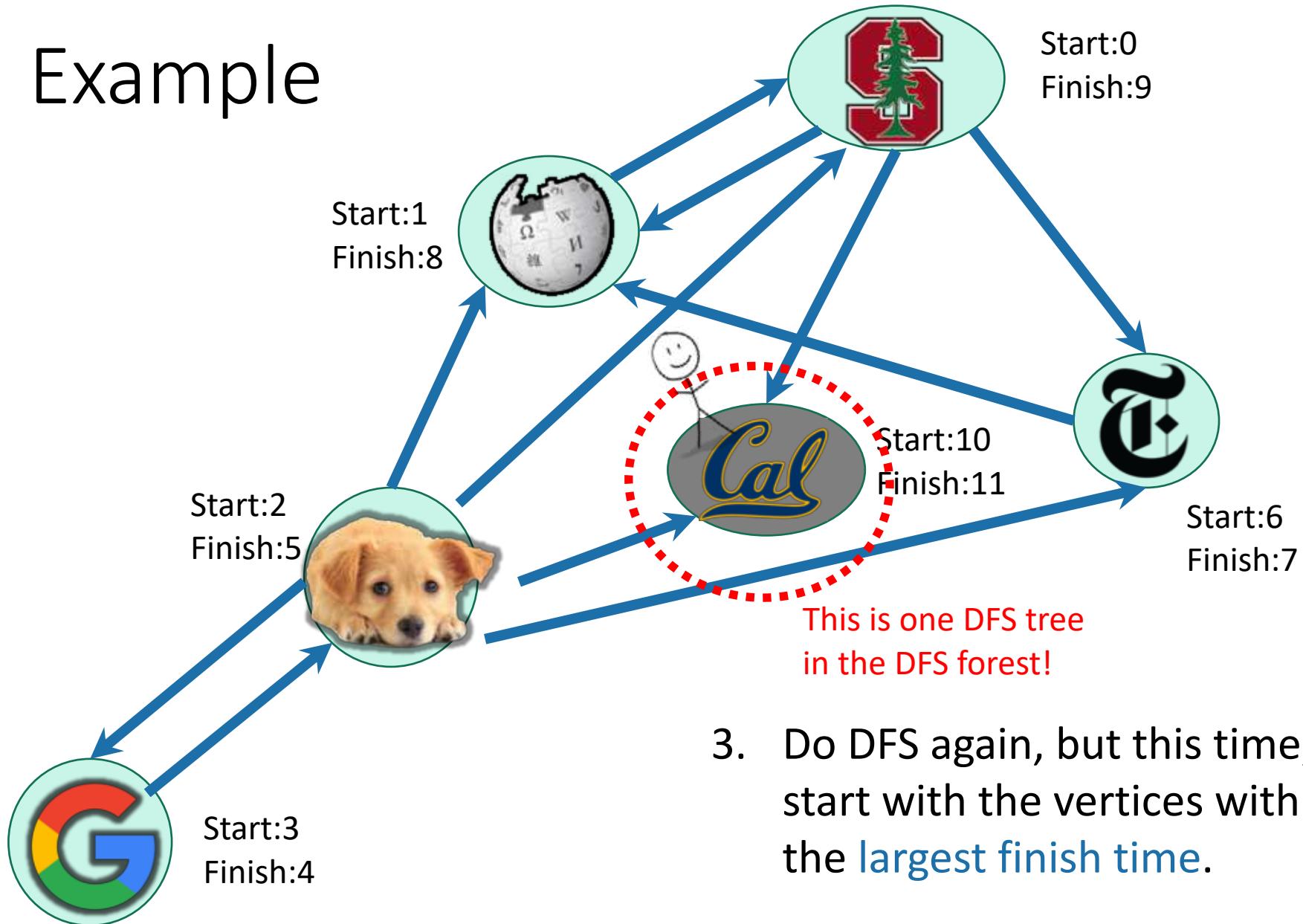
Example



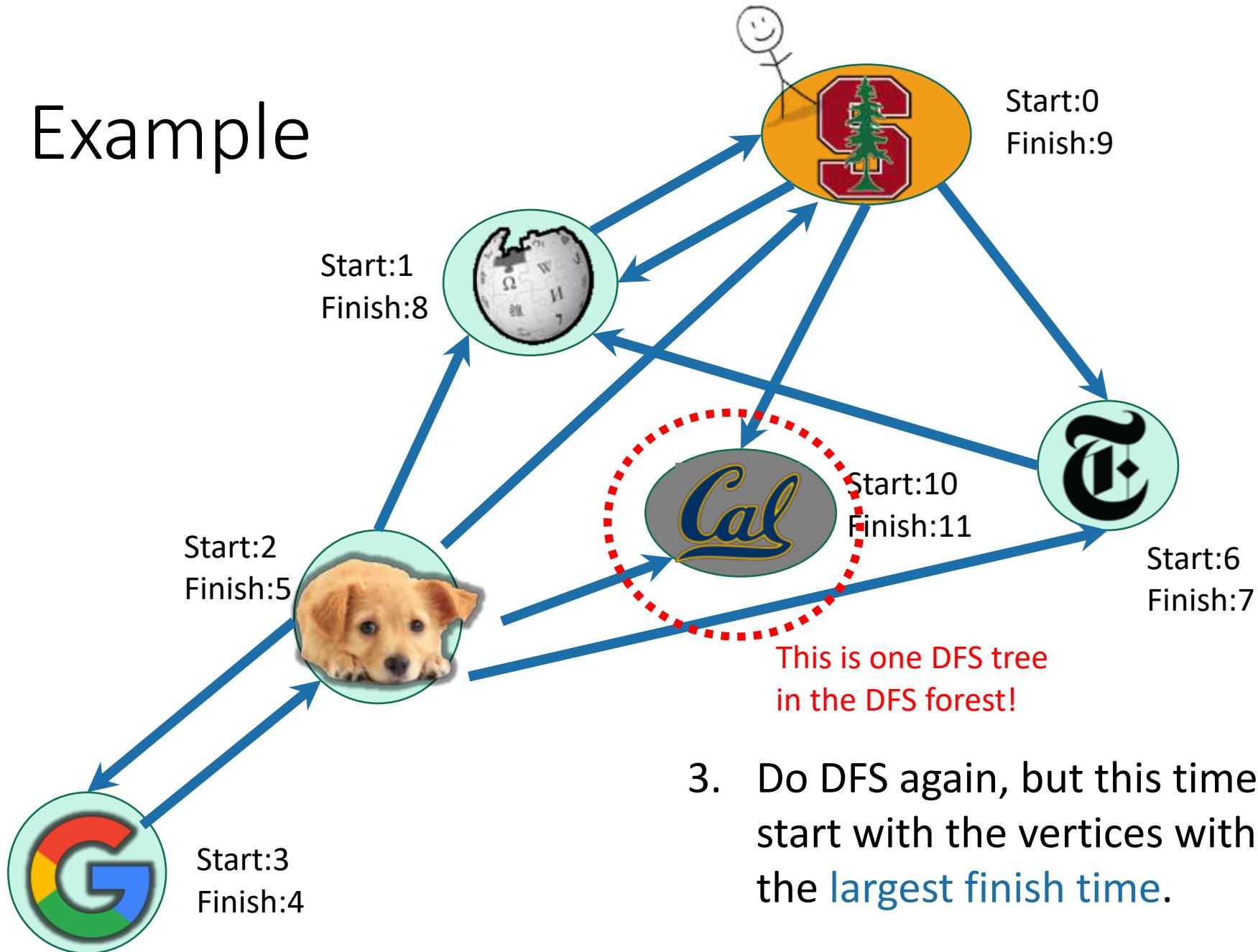
Example



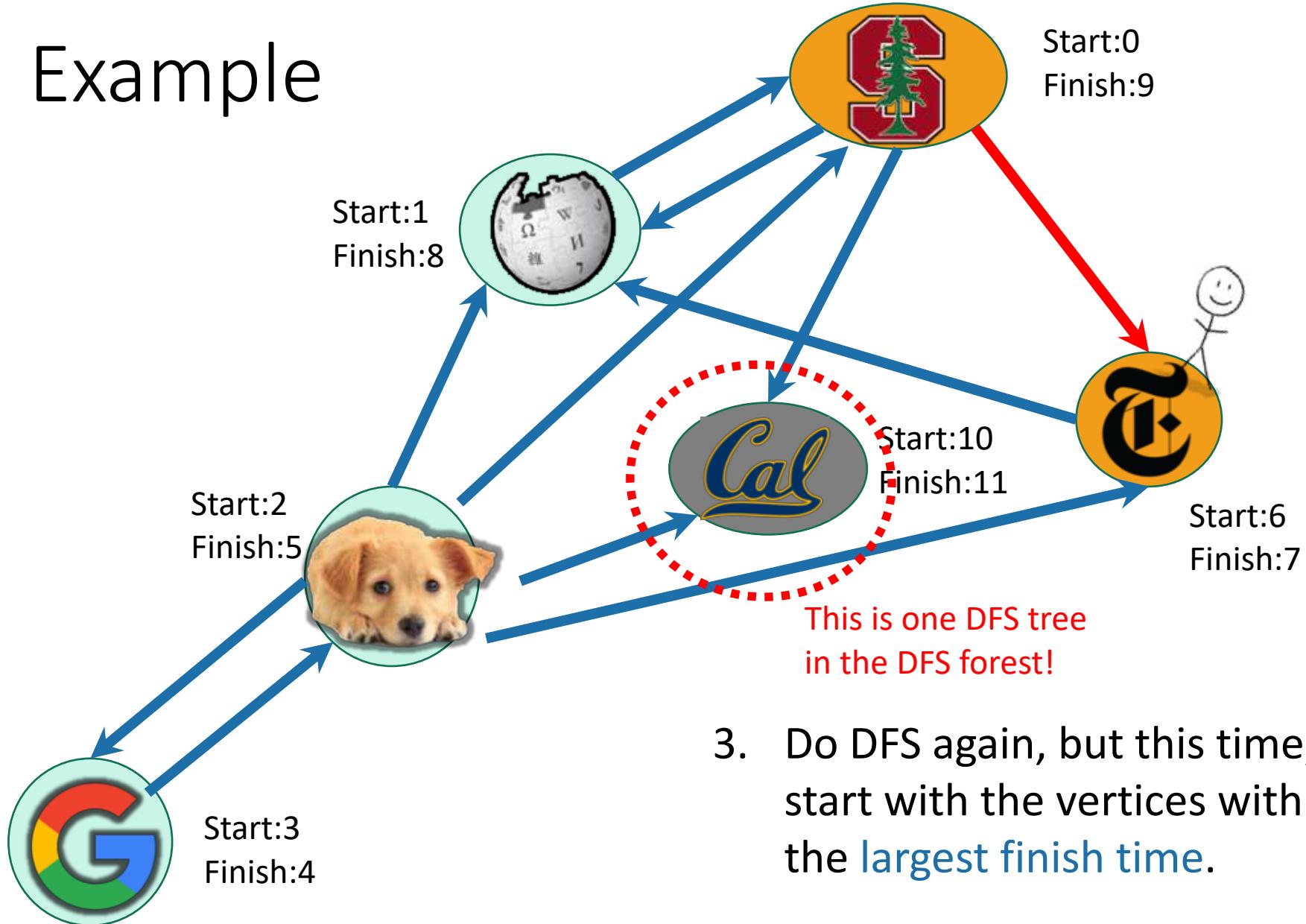
Example



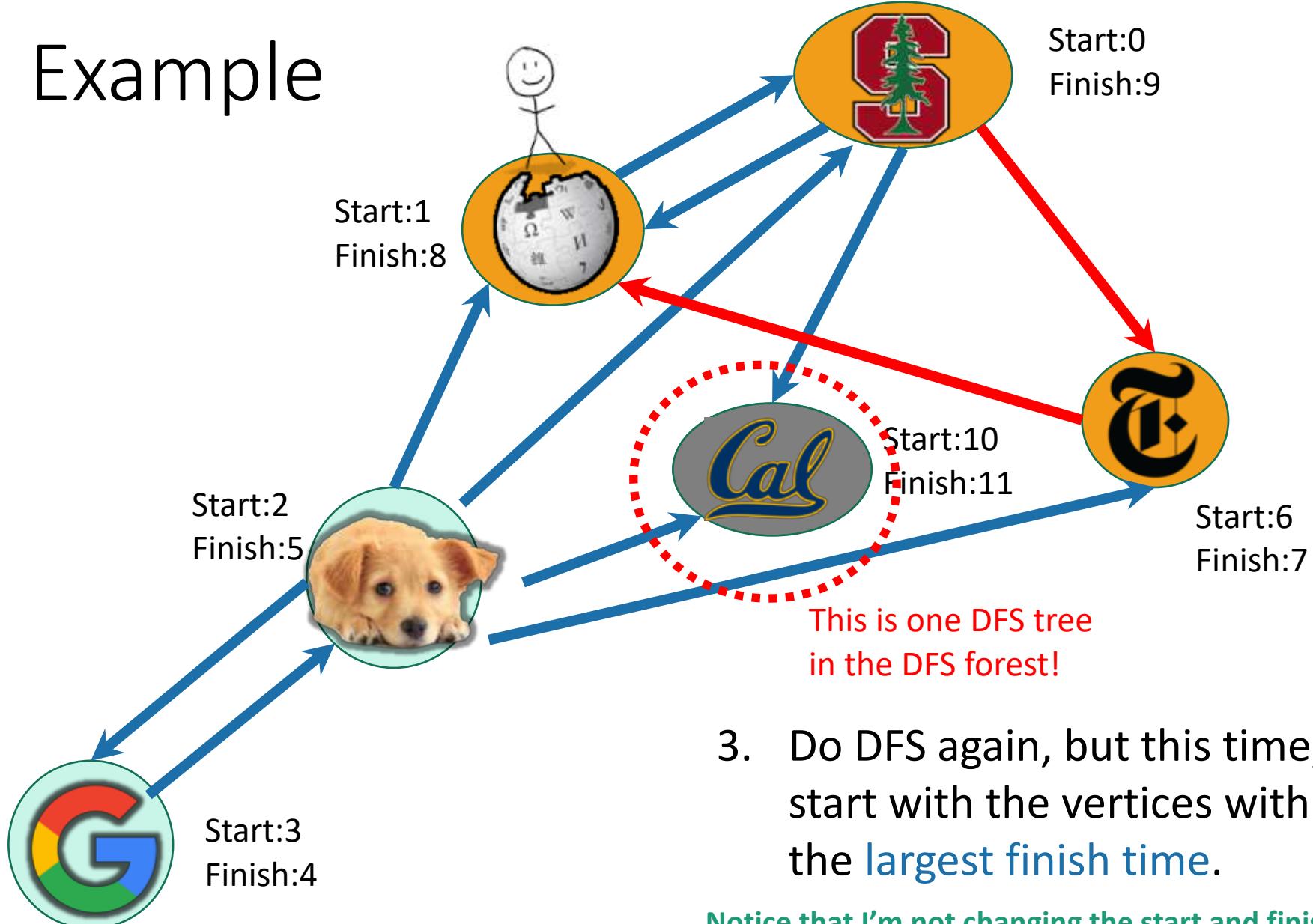
Example



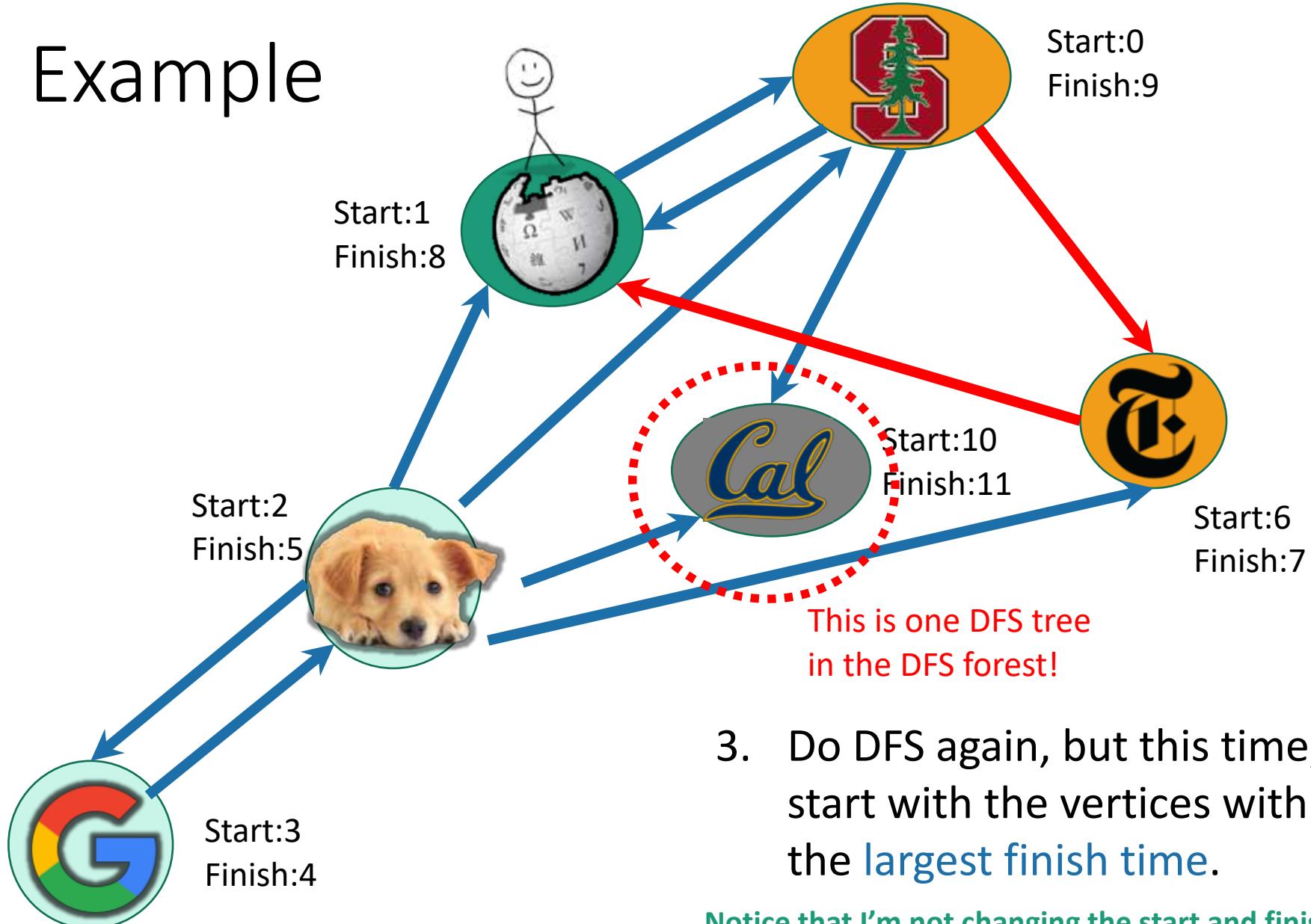
Example



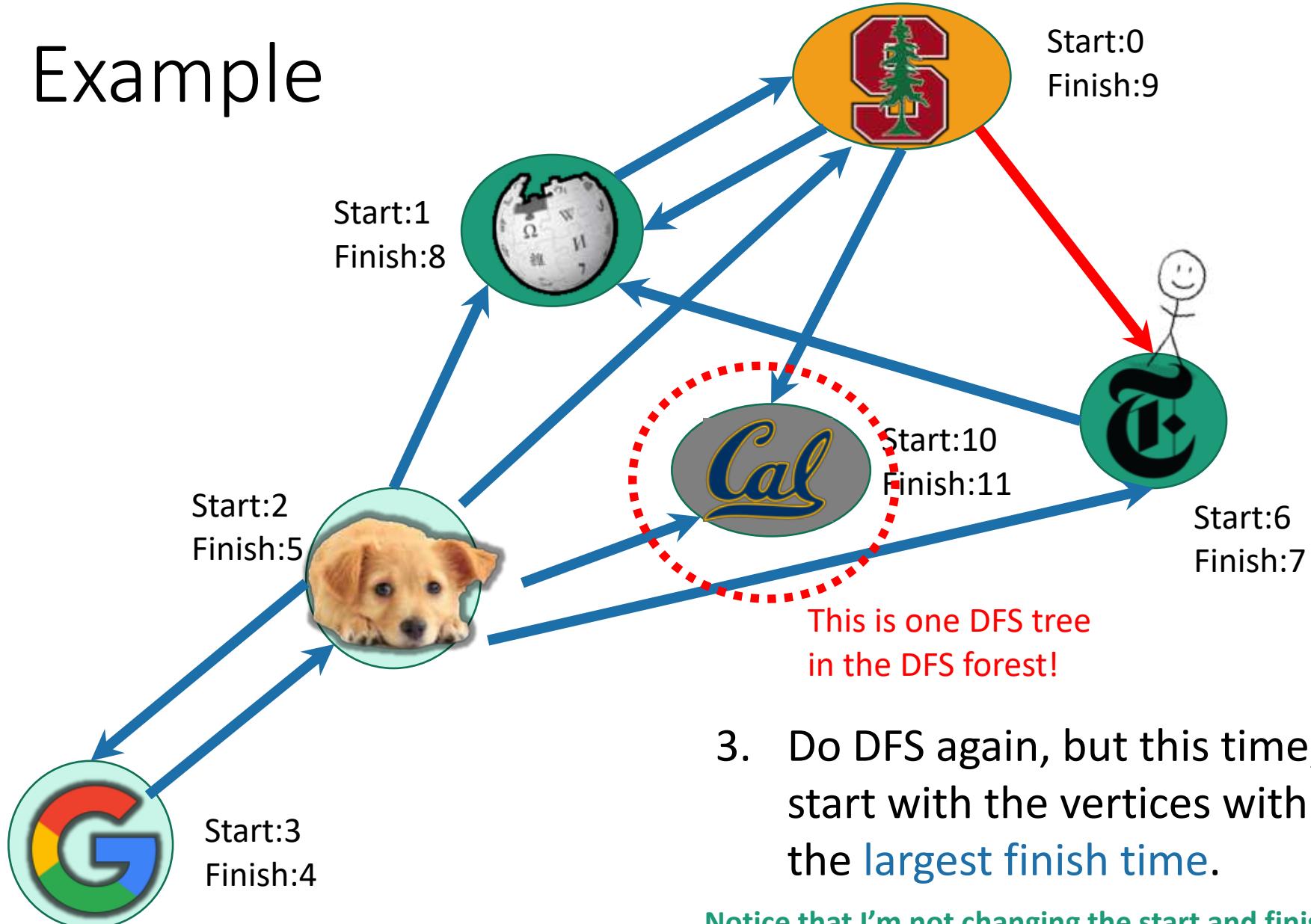
Example



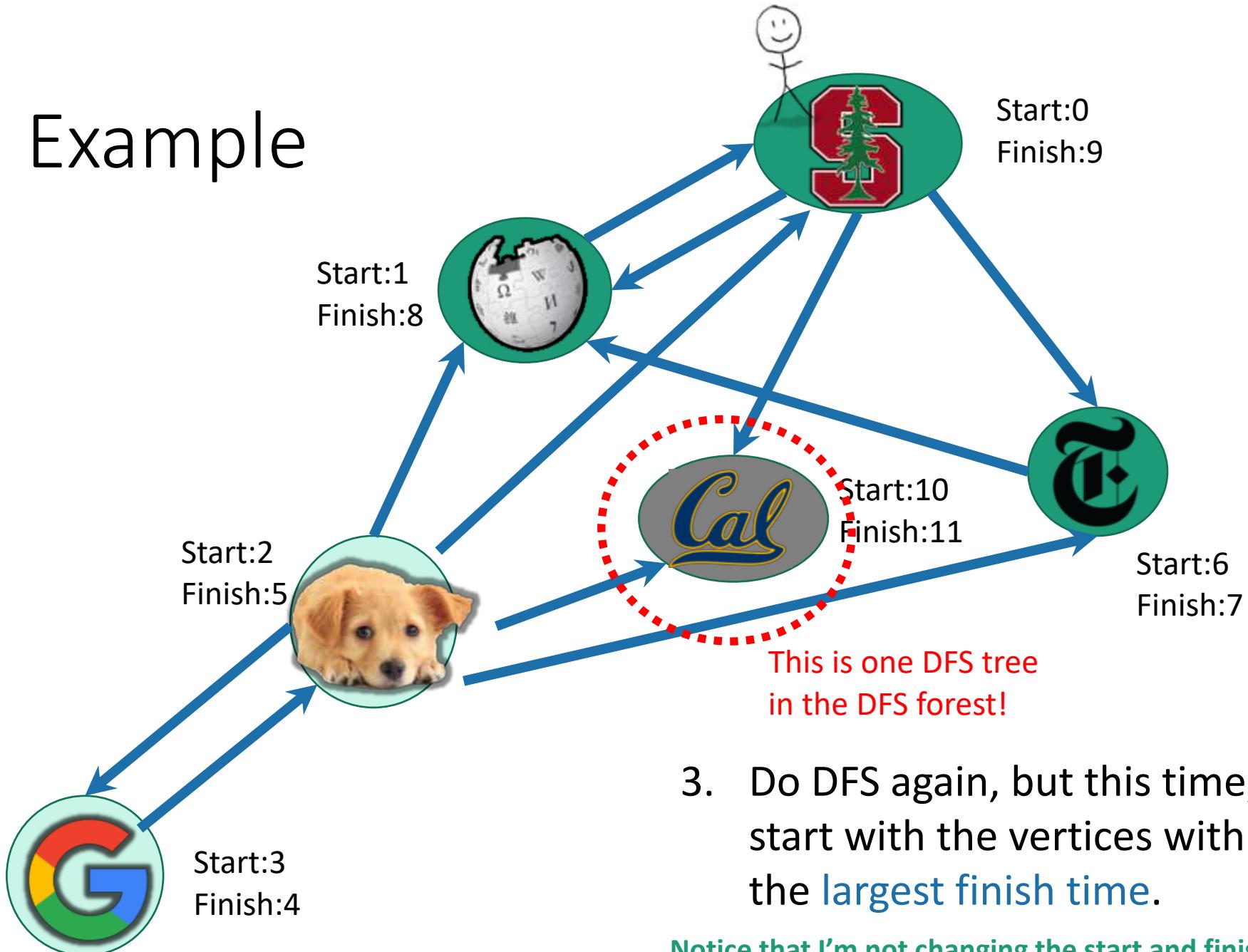
Example



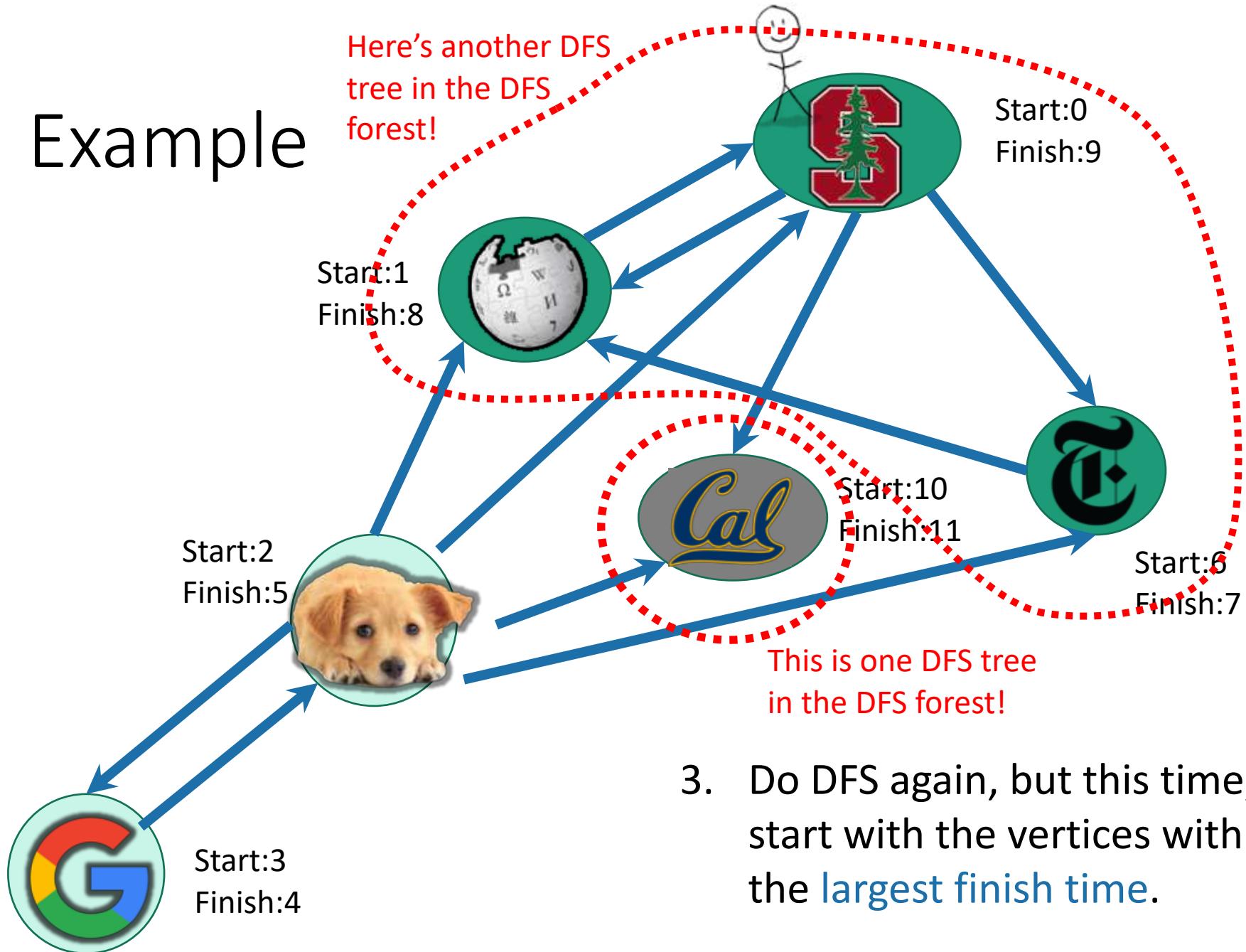
Example



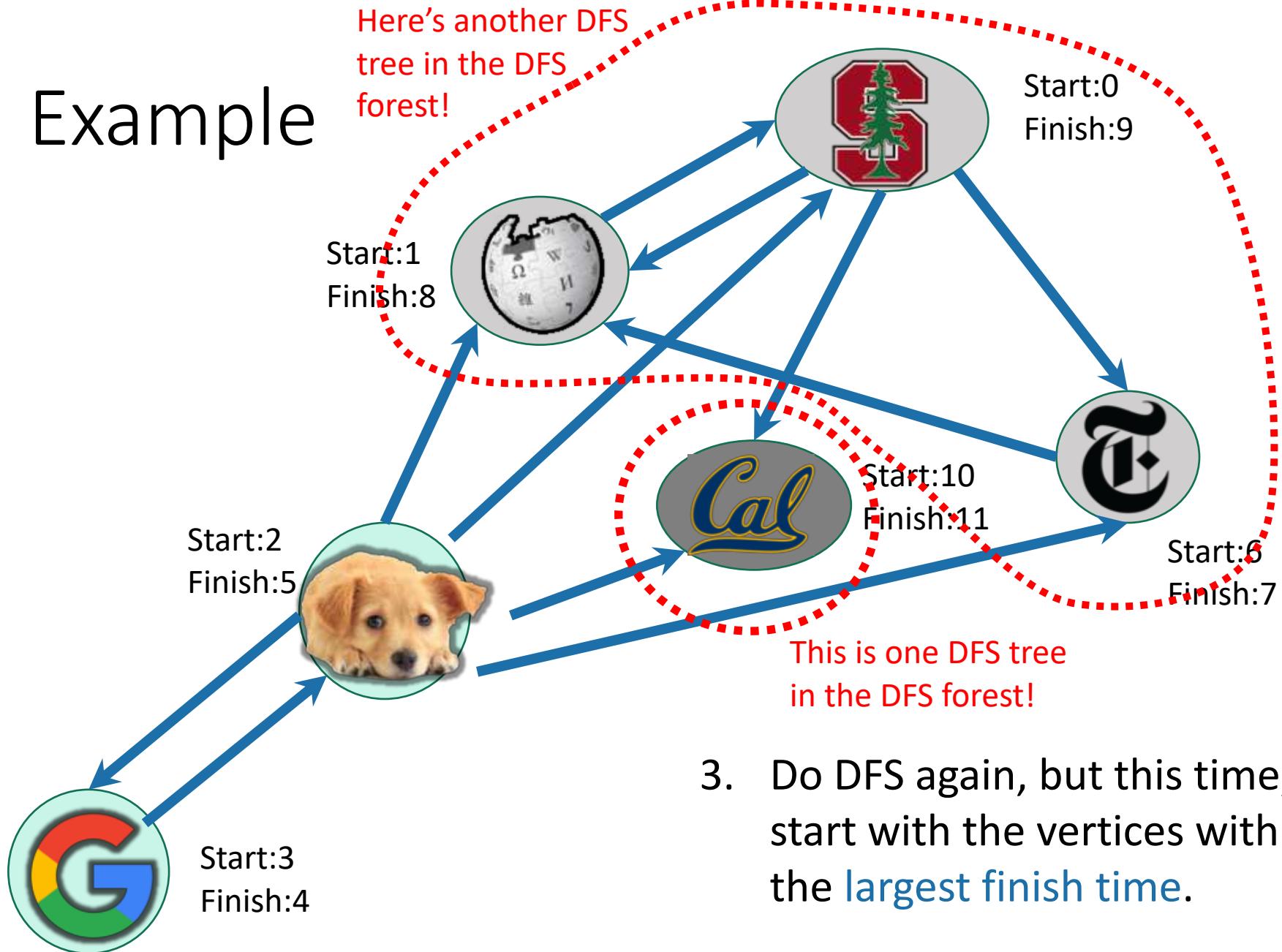
Example



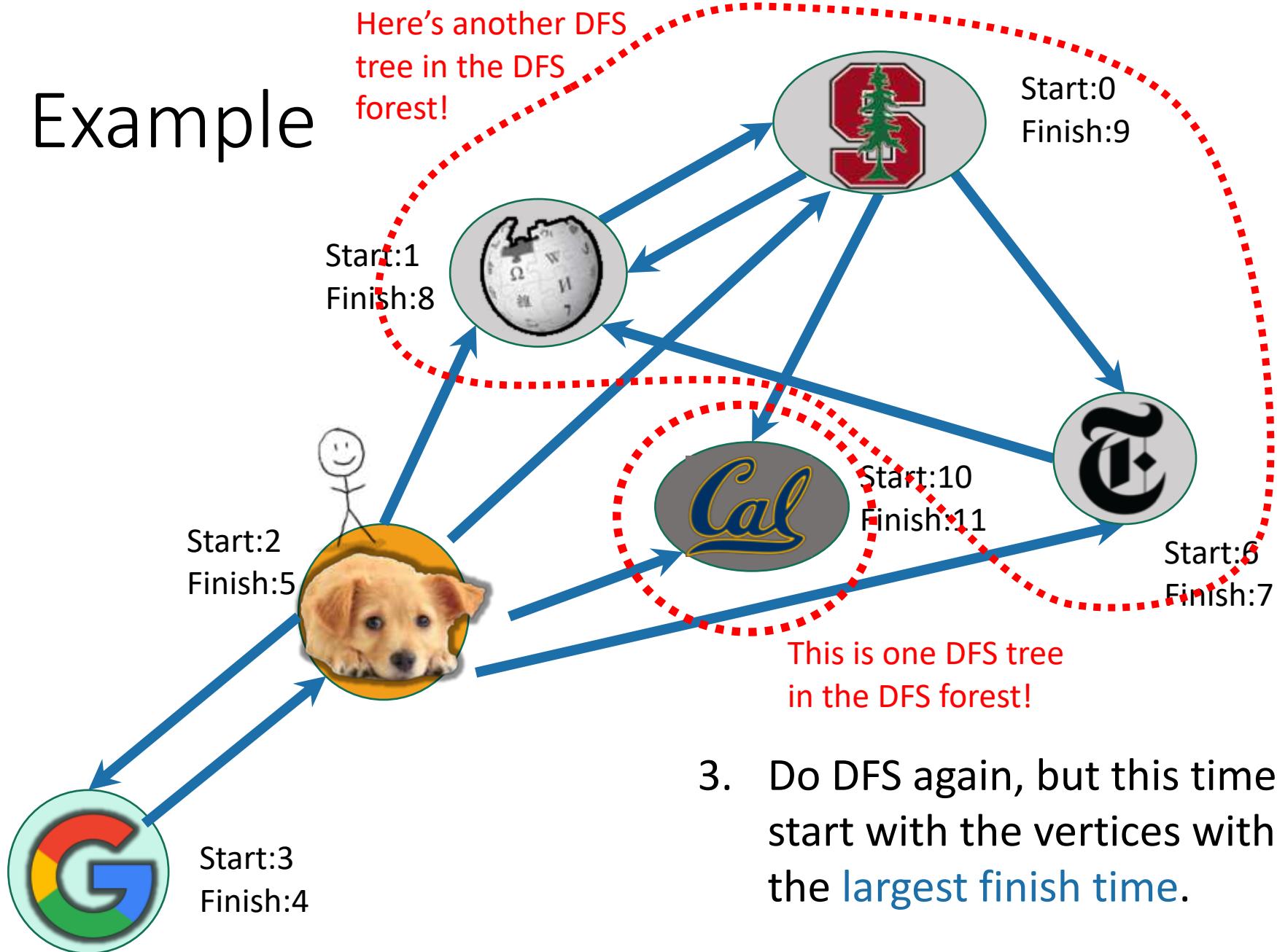
Example



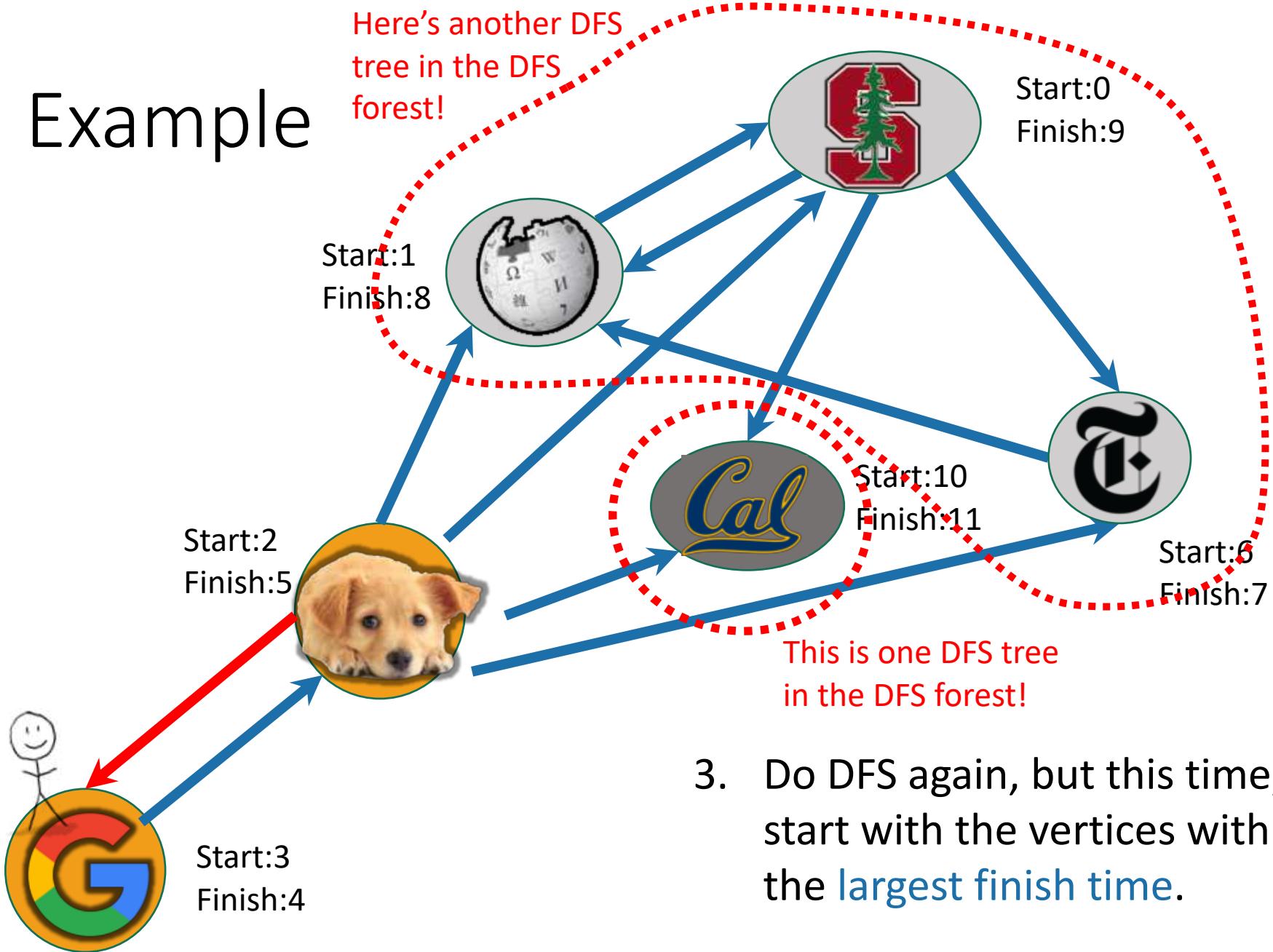
Example



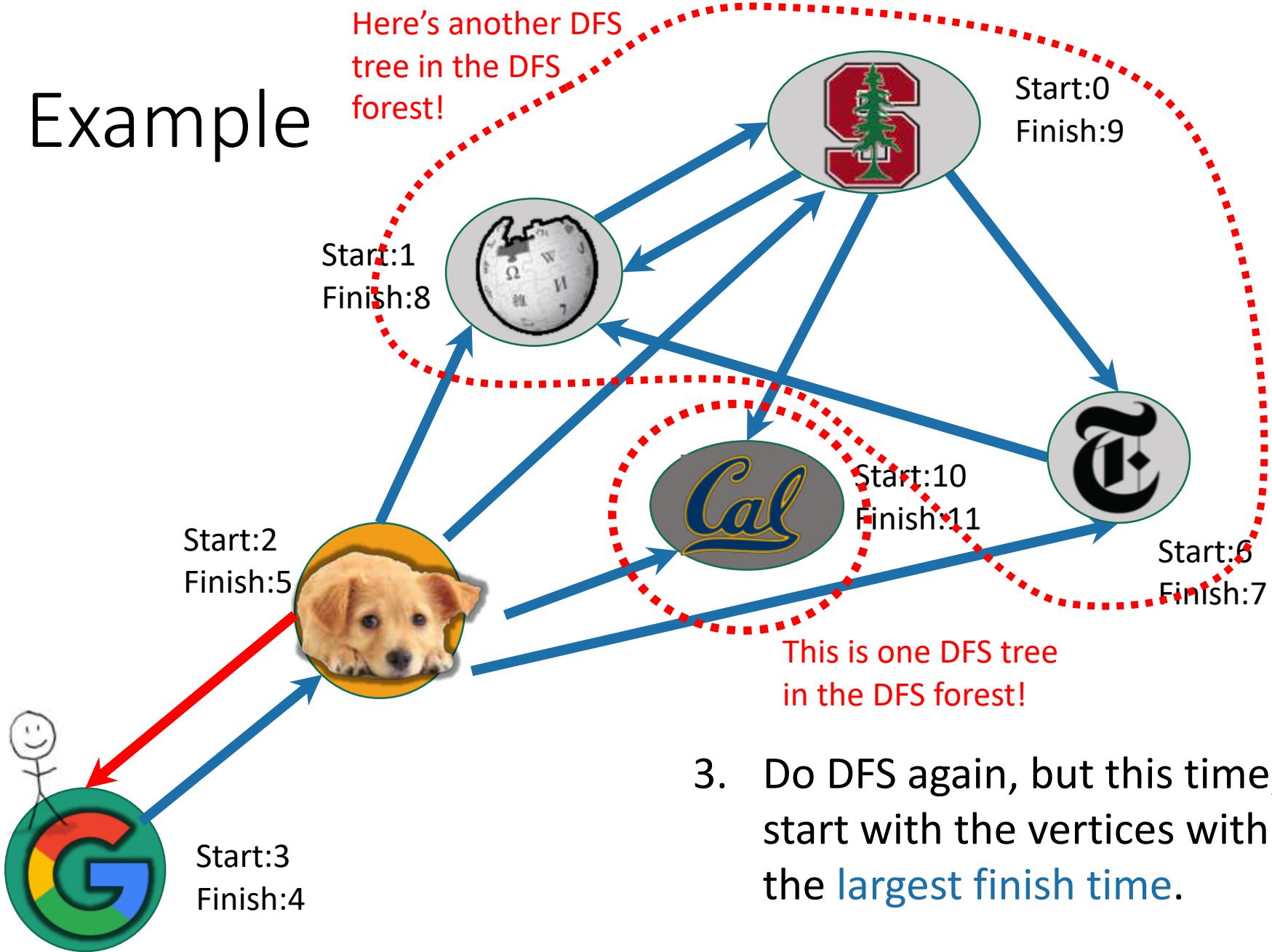
Example



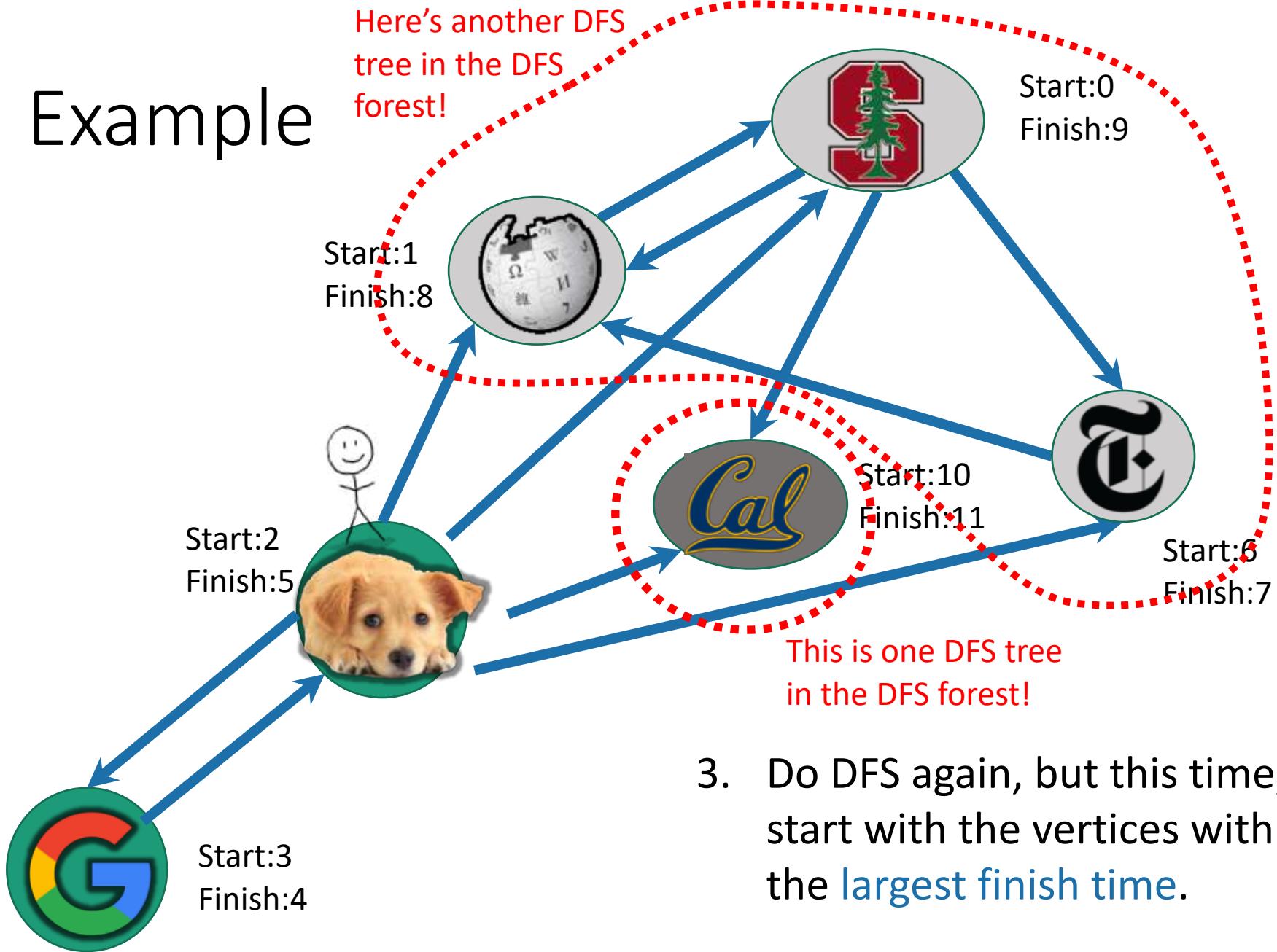
Example



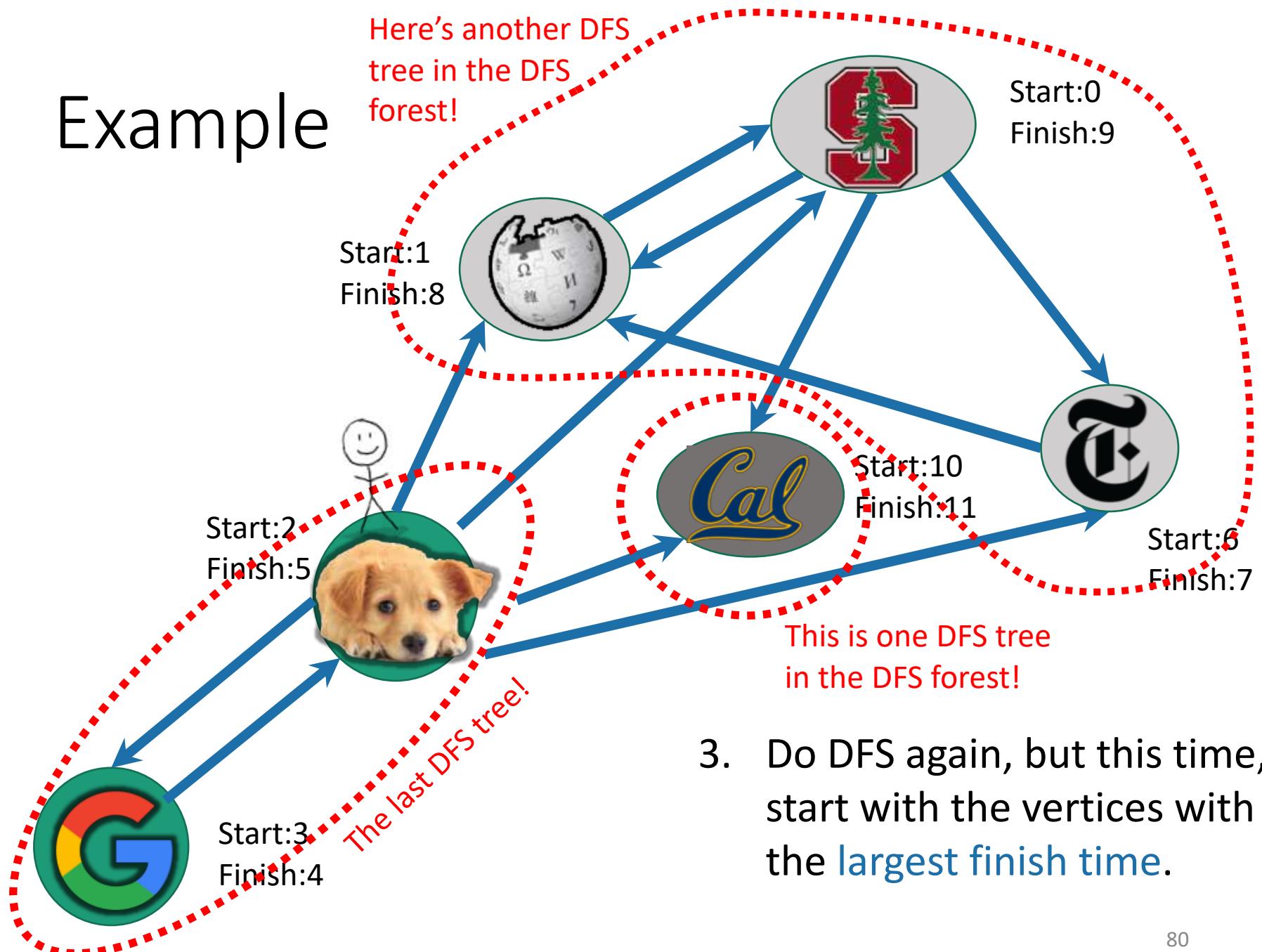
Example



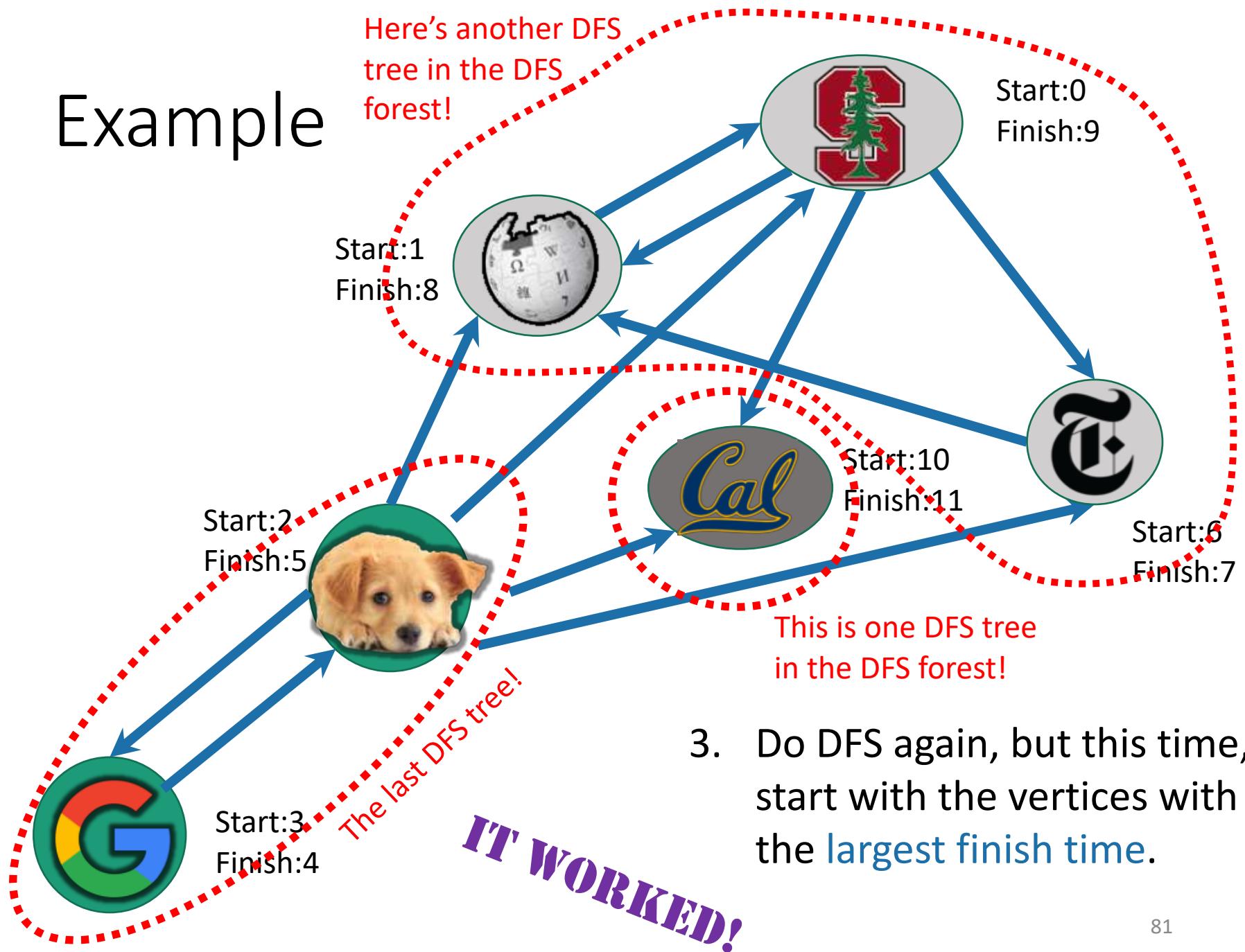
Example



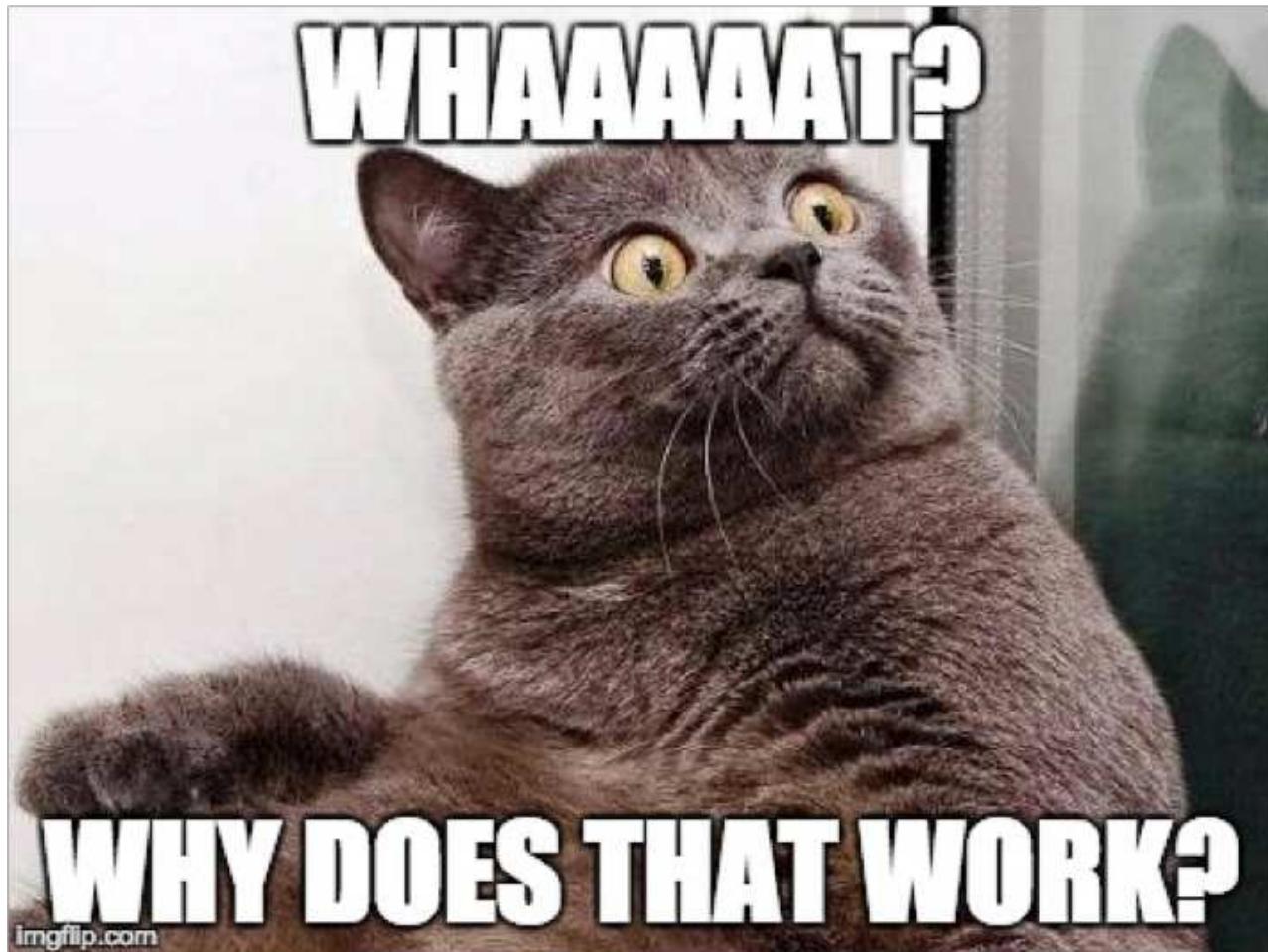
Example



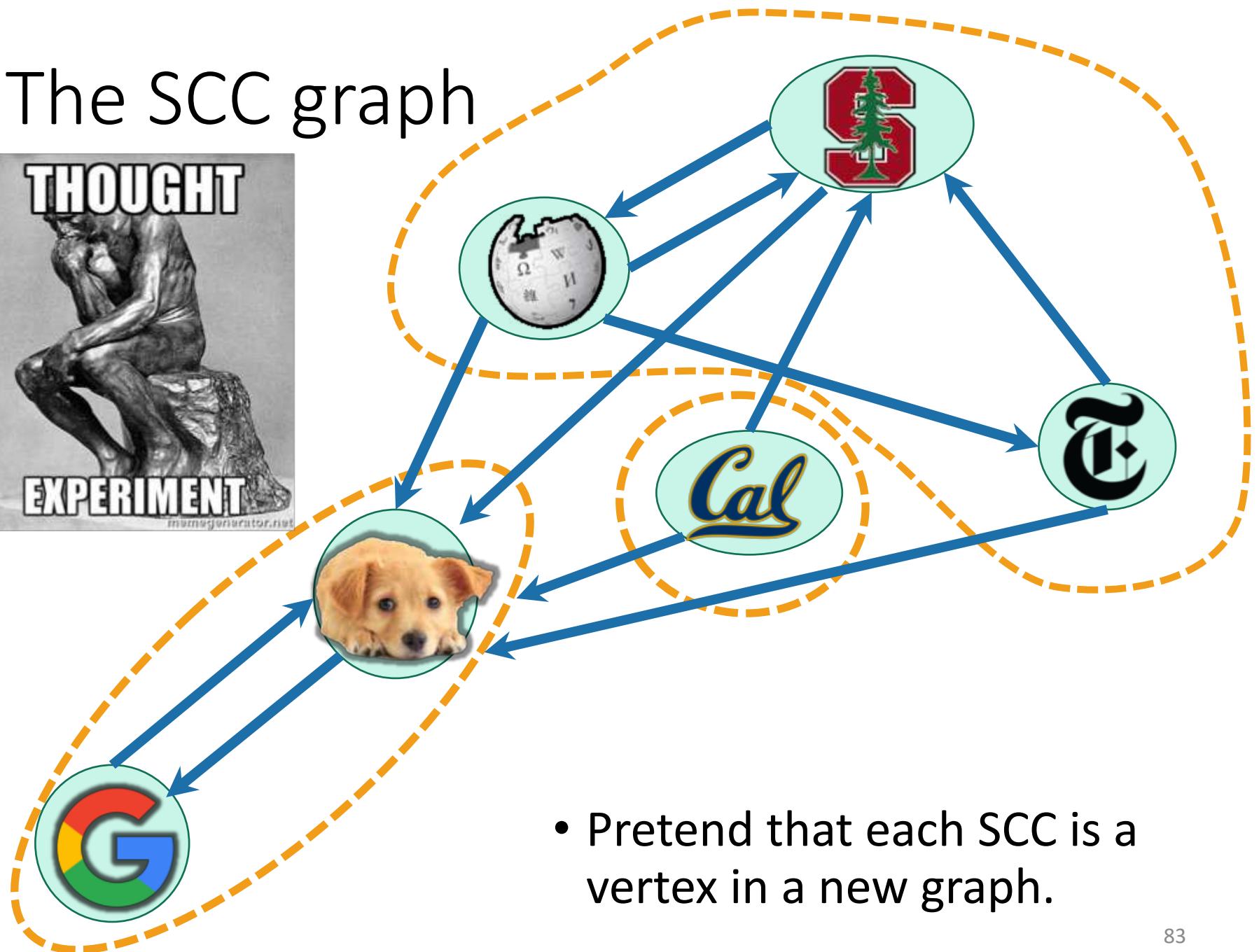
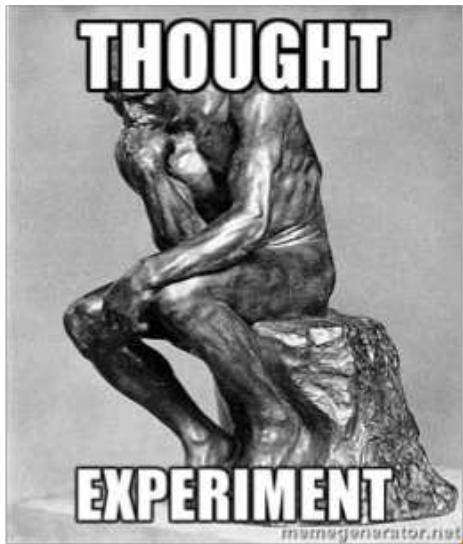
Example



One question



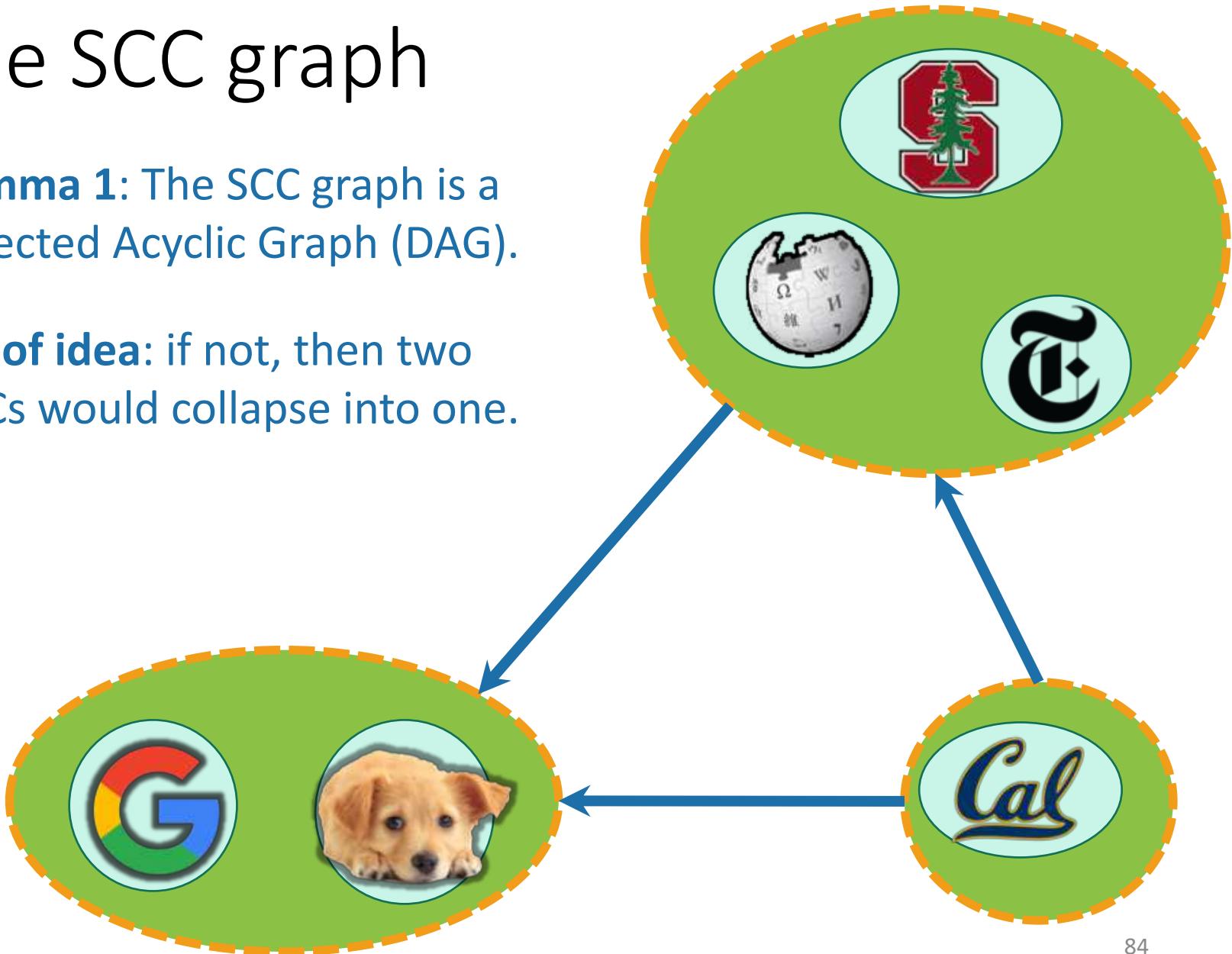
The SCC graph



The SCC graph

Lemma 1: The SCC graph is a Directed Acyclic Graph (DAG).

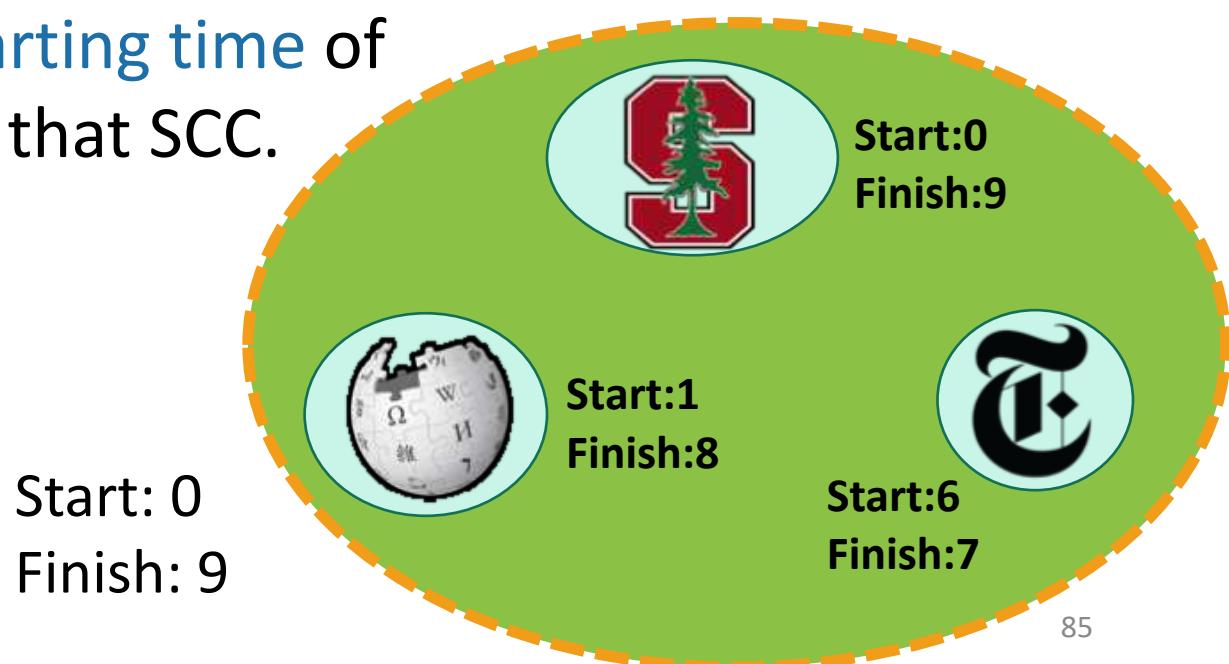
Proof idea: if not, then two SCCs would collapse into one.



Starting and finishing times in a SCC

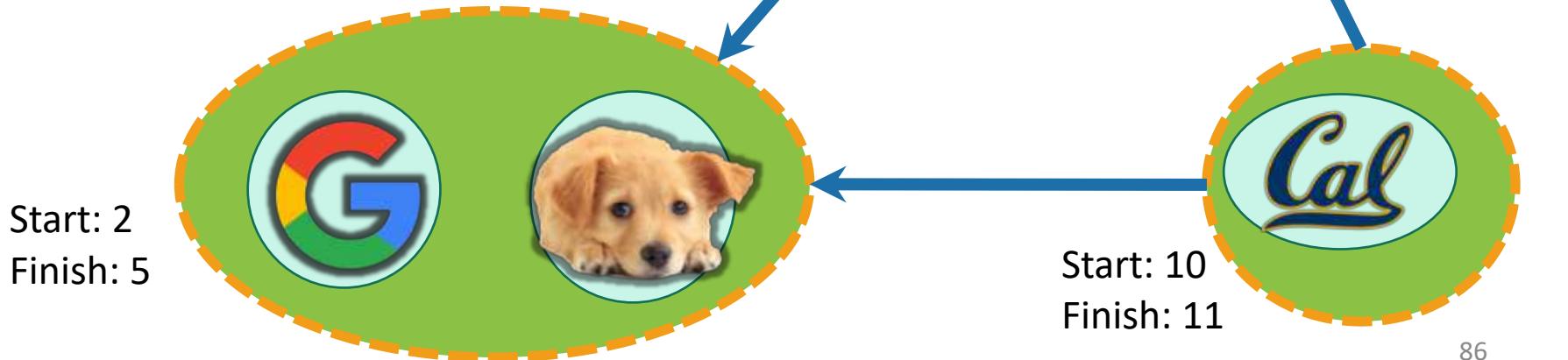
Definitions:

- The **finishing time** of a SCC is the **largest finishing time** of any element of that SCC.
- The **starting time** of a SCC is the **smallest starting time** of any element of that SCC.



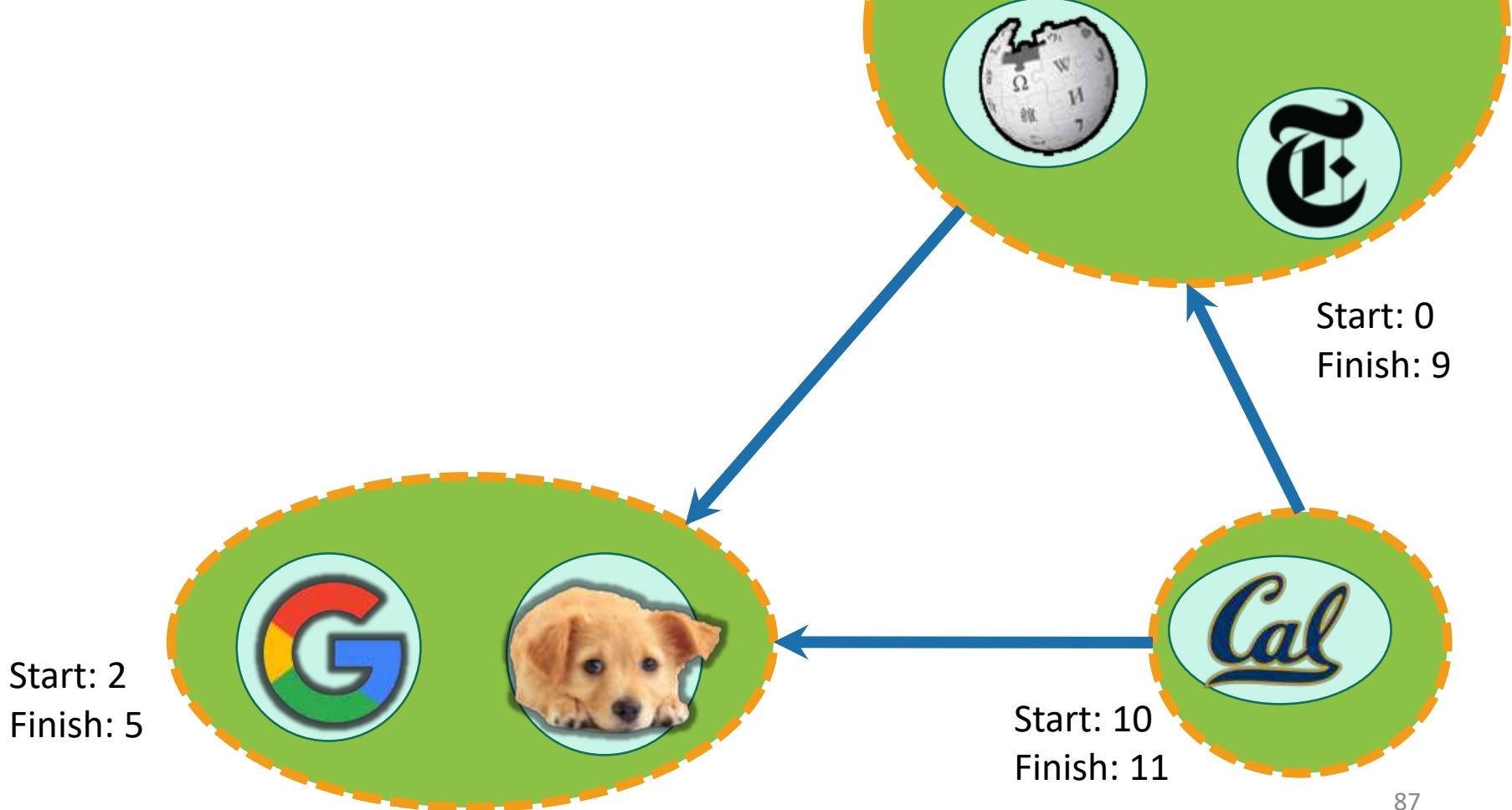
Our SCC DAG with start and finish times

- Last time we saw that Finishing times allowed us to **topologically sort** of the vertices.
- Notice that works in this example too...



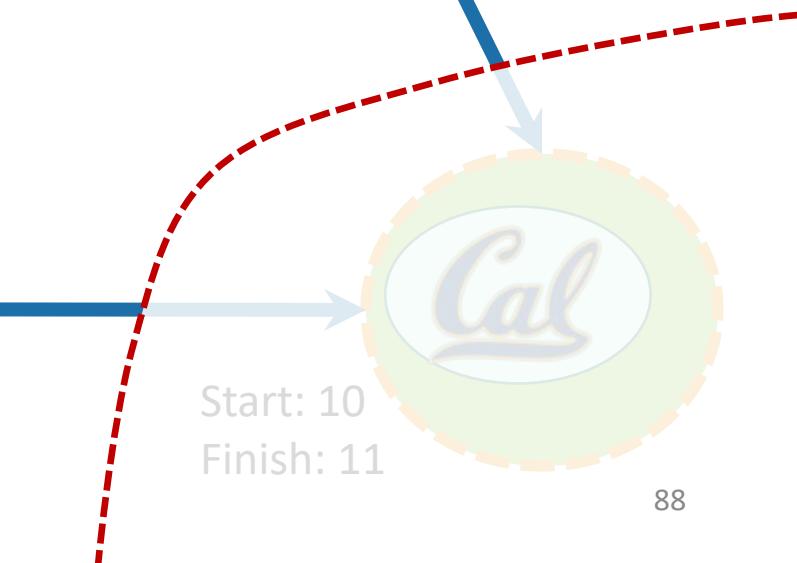
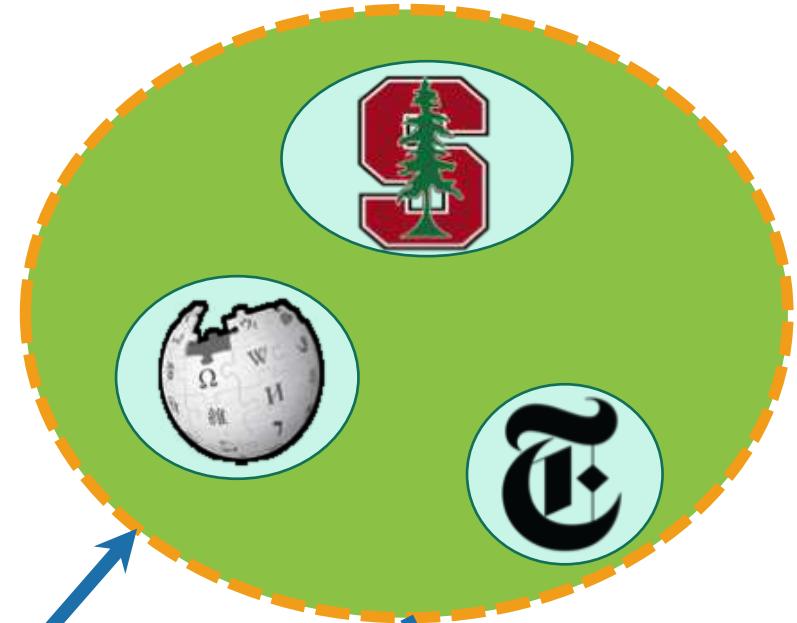
Main idea

- Let's reverse the edges.



Main idea

- Let's reverse the edges.
- Now, the SCC with the largest finish time has no edges going out.
 - If it did have edges going out, then it wouldn't be a good thing to choose first in a topological ordering!
- If I run DFS there, I'll find exactly that component.
- Remove and repeat.



Let's make this idea formal.

Back the the parentheses theorem

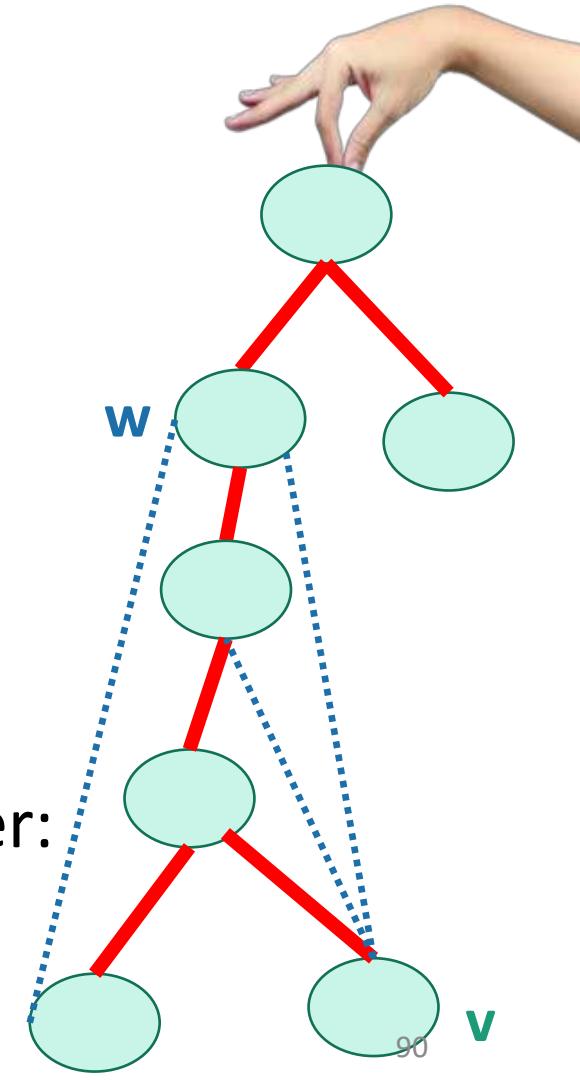
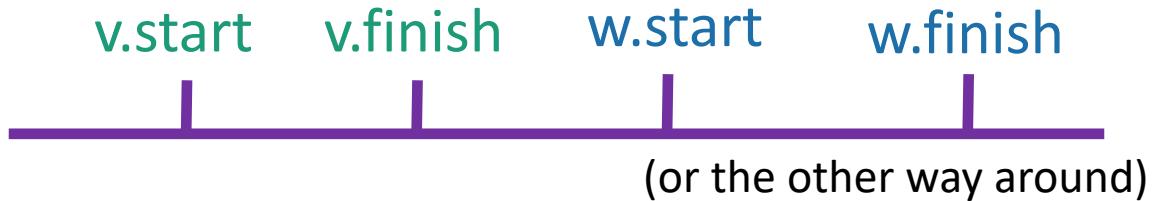
- If v is a descendent of w in this tree:



- If w is a descendent of v in this tree:

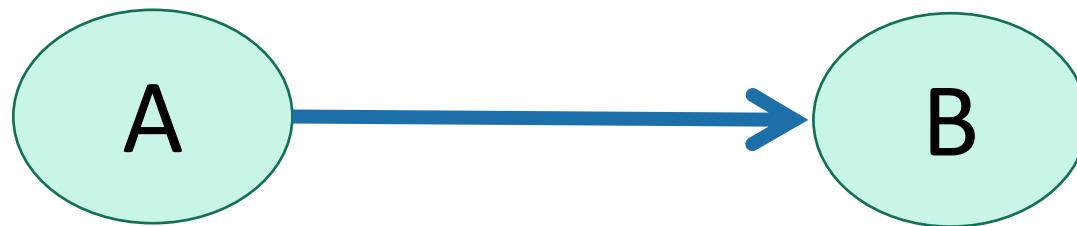


- If neither are descendants of each other:



As we saw last time...

Claim: In a DAG, we'll always have:

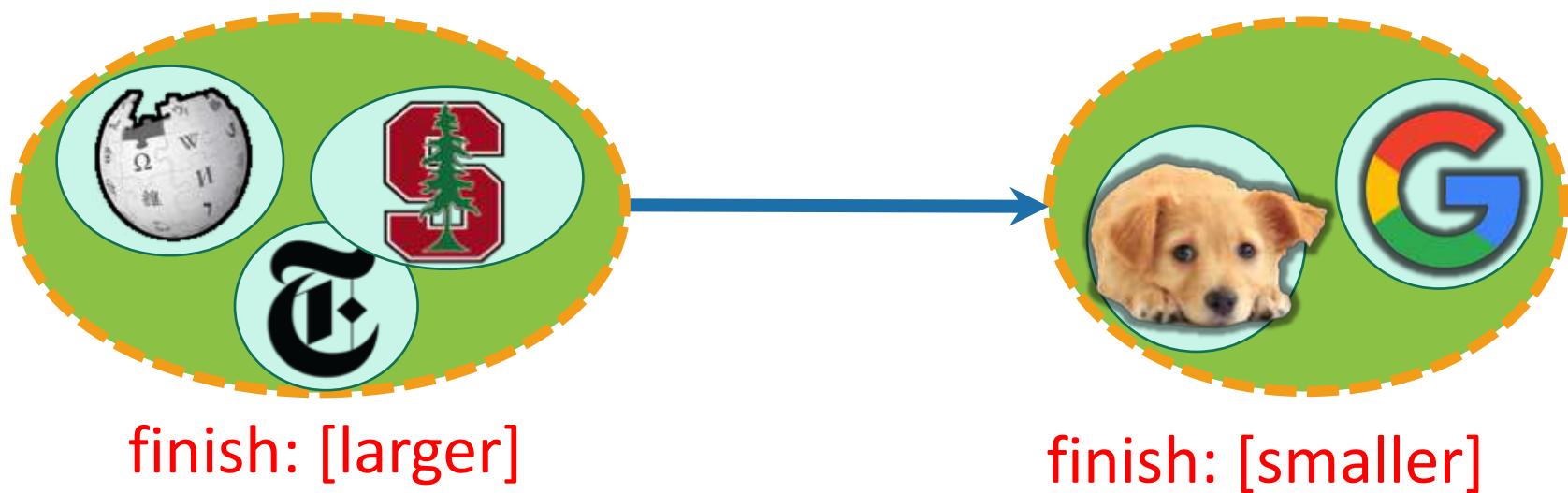


finish: [larger]

finish: [smaller]

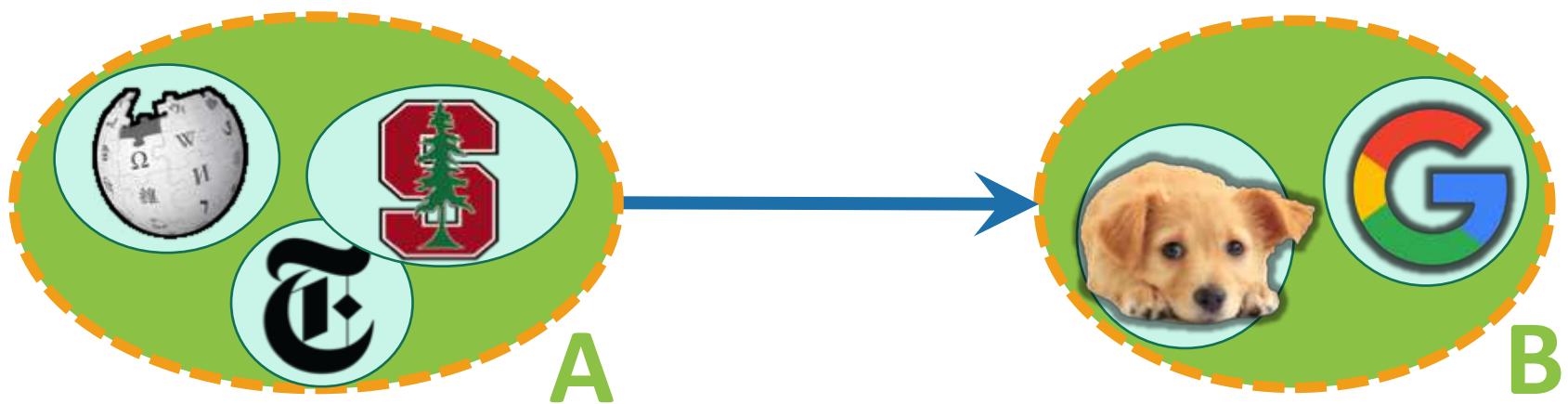
Same thing, in the SCC DAG.

- **Claim:** we'll always have



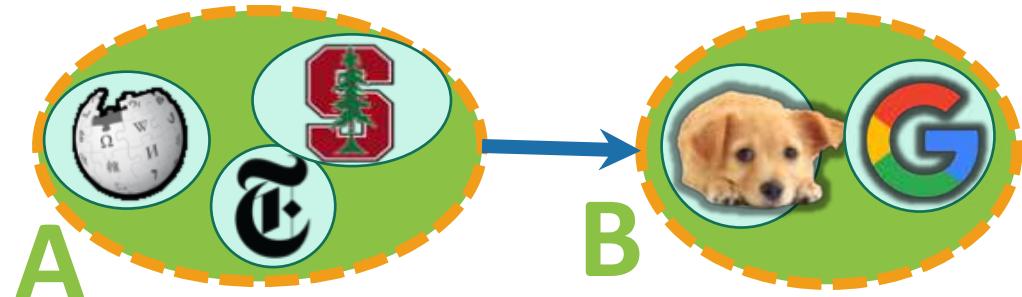
Let's call it Lemma 2

- If there is an edge like this:



- Then $A.\text{finish} > B.\text{finish}$.

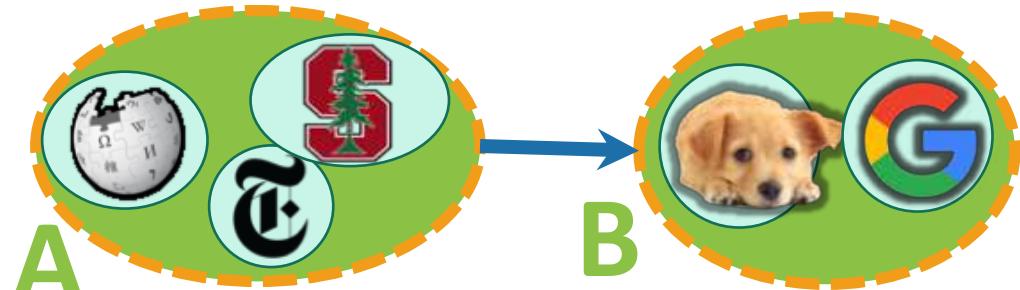
Proof idea



Want to show $A.\text{finish} > B.\text{finish}$.

- **Two cases:**
 - We reached **A** before **B** in our first DFS.
 - We reached **B** before **A** in our first DFS.

Proof idea



Want to show $A.\text{finish} > B.\text{finish}$.

- **Case 1:** We reached **A** before **B** in our first DFS.

- Say that:

- **x** has the largest finish time in **A**;
- **y** has the largest finish in **B**;
- **z** was discovered first in **A**;

- Then:

- Reach **A** before **B**
- \Rightarrow we will discover **y** via **z**
- \Rightarrow **y** is a descendant of **z** in the DFS forest.

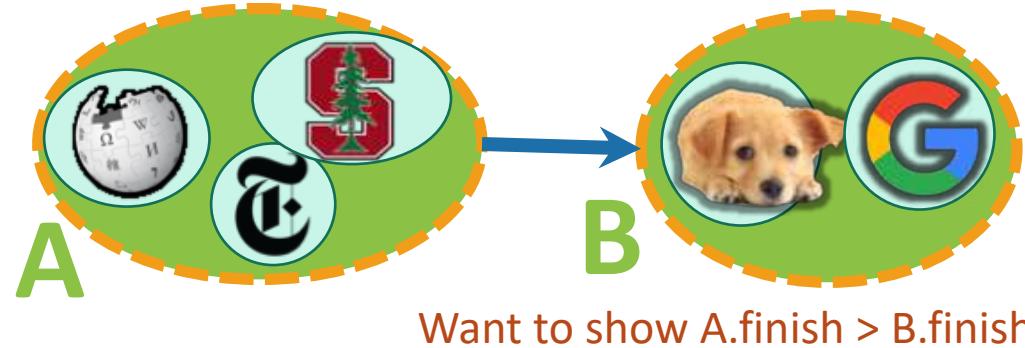
So $A.\text{finish} = x.\text{finish}$
 $B.\text{finish} = y.\text{finish}$
 $x.\text{finish} \geq z.\text{finish}$

aka,
A.finish > B.finish

- Then



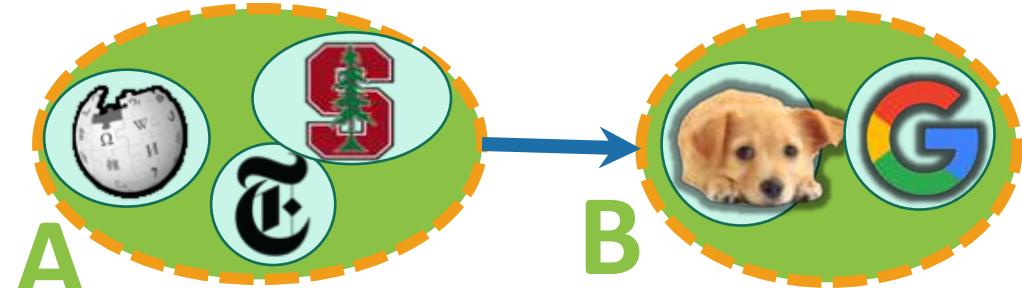
Proof idea



- **Case 2:** We reached **B** before **A** in our first DFS.
- There are no paths from B to A
 - because the SCC graph has no cycles
- So we completely finish exploring B and never reach A.
- A is explored later after we restart DFS.

aka,
A.finish > B.finish

Proof idea



Want to show $A.\text{finish} > B.\text{finish}$.

- **Two cases:**
 - We reached **A** before **B** in our first DFS.
 - We reached **B** before **A** in our first DFS.
- In either case:

A.finish > B.finish

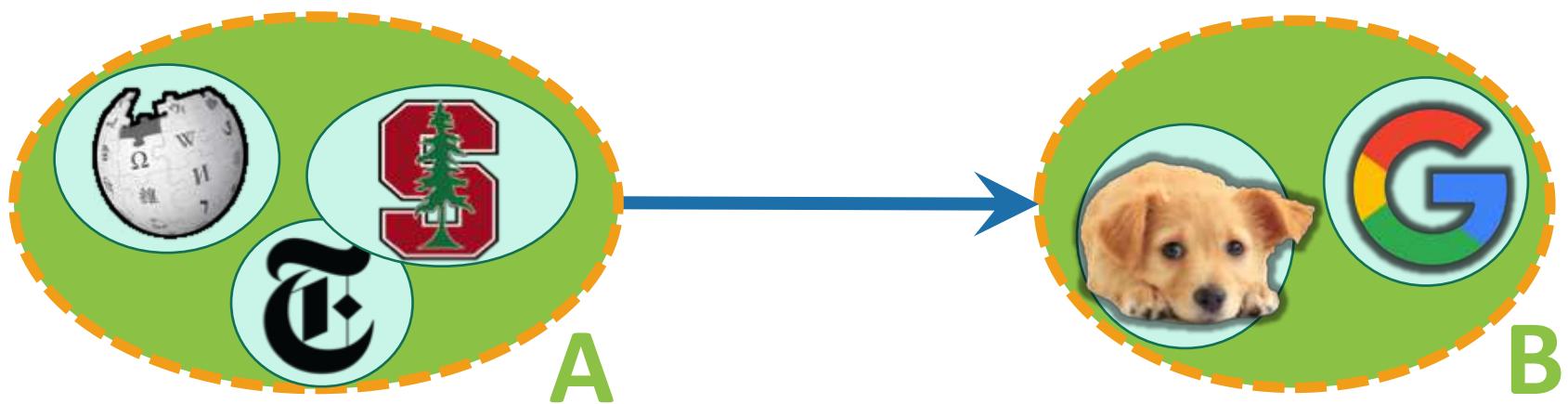
which is what we wanted to show.



Notice: this is exactly the same two-case argument that we did last time for topological sorting, just with the SCC DAG!

This establishes:
Lemma 2

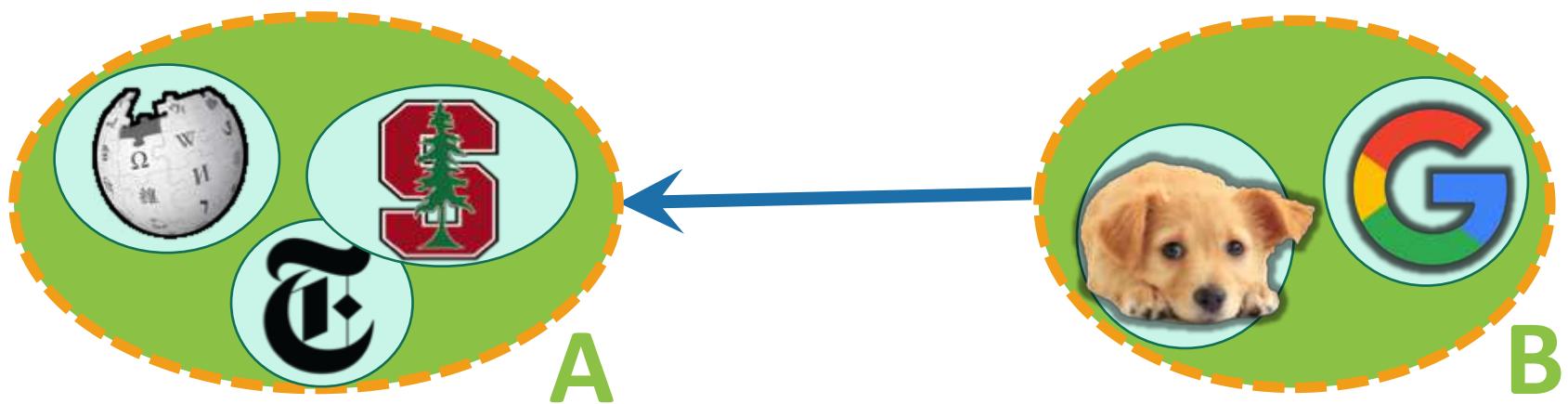
- If there is an edge like this:



- Then $A.\text{finish} > B.\text{finish}$.

This establishes:
Corollary 1

- If there is an edge like this in the **reversed graph**:

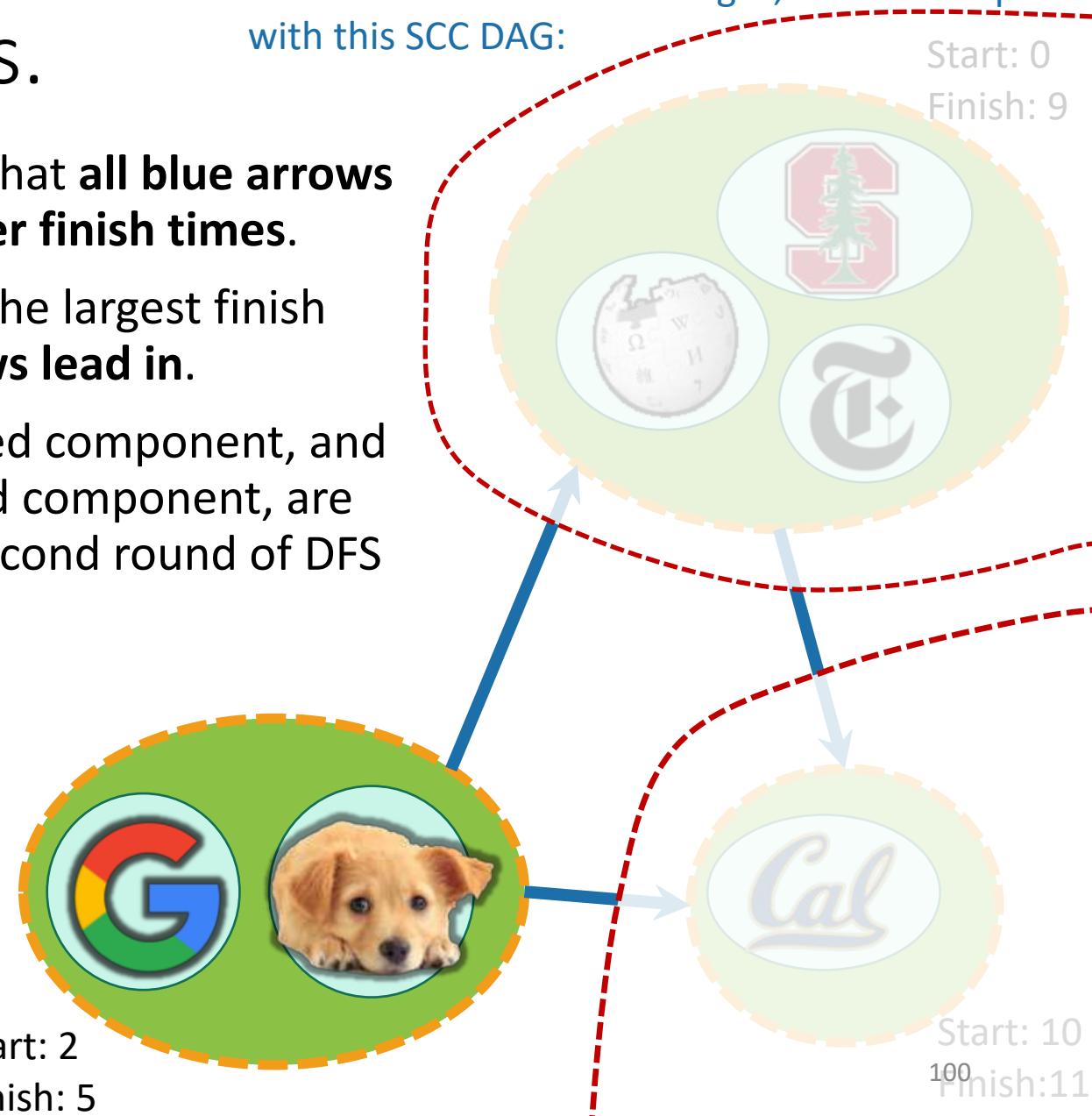


- Then $A.\text{finish} > B.\text{finish}$.

Now we see why this finds SCCs.

- The Corollary says that **all blue arrows point towards larger finish times**.
 - So if we start with the largest finish time, **all blue arrows lead in**.
 - Thus, that connected component, and only that connected component, are reachable by the second round of DFS
-
- Now, we've deleted that first component.
 - The next one has the **next biggest finishing time**.
 - So **all remaining blue arrows lead in**.
 - **Repeat.**

Remember that after the first round of DFS, and after we reversed all the edges, we ended up with this SCC DAG:



Formally, we prove it by induction

- **Theorem:** The algorithm we saw before will correctly identify strongly connected components.
- **Inductive hypothesis:**
 - The first t trees found in the second (reversed) DFS forest are the t SCCs with the largest finish times.
 - Moreover, what's left unvisited after these t trees have been explored is a DAG on the un-found SCCs.
- **Base case: ($t=0$)**
 - The first 0 trees found in the reversed DFS forest are the 0 SCCs with the largest finish times. (**TRUE**)
 - Moreover, what's left unvisited after 0 trees have been explored is a DAG on all the SCCs. (**TRUE by Lemma 1.**)

Inductive step [drawing on board to supplement]

- Assume by induction that the first t trees are the last-finishing SCCs, and the remaining SCCs form a DAG.
- Consider the $(t+1)^{\text{st}}$ tree produced, suppose the root is $\textcolor{brown}{x}$.
- Suppose that $\textcolor{brown}{x}$ lives in the SCC $\textcolor{green}{A}$.
- Then $\textcolor{green}{A}.\text{finish} > \textcolor{orange}{B}.\text{finish}$ for all remaining SCCs $\textcolor{orange}{B}$.
 - This is because we chose $\textcolor{brown}{x}$ to have the largest finish time.
- Then there are no edges leaving $\textcolor{green}{A}$ in the remaining SCC DAG.
 - This follows from the Corollary.
- Then DFS started at $\textcolor{brown}{x}$ recovers exactly $\textcolor{green}{A}$.
 - It doesn't recover any more since nothing else is reachable.
 - It doesn't recover any less since A is strongly connected.
 - (Notice that we are using that A is still strongly connected when we reverse all the edges).
- So the $(t+1)^{\text{st}}$ tree is the SCC with the $(t+1)^{\text{st}}$ biggest finish time.
(Also the remaining SCCs still form a DAG, since removing a vertex won't create cycles).

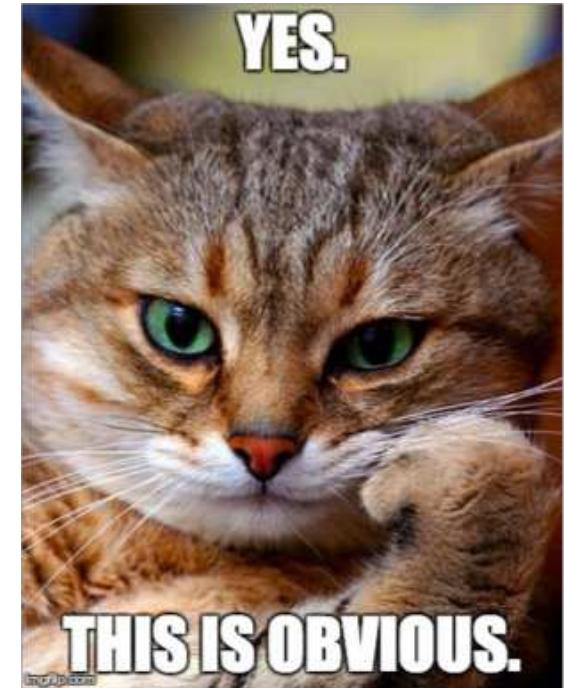
Formally, we prove it by induction

- **Theorem:** The algorithm we saw before will correctly identify strongly connected components.
- **Inductive hypothesis:**
 - The first t trees found in the second (reversed) DFS forest are the t SCCs with the largest finish times.
 - Moreover, what's left unvisited after these t trees have been explored is a DAG on the un-found SCCs.
- **Base case:** [done]
- **Inductive step:** [done]
- **Conclusion:** The second (reversed) DFS forest contains all the SCCs as its trees!
 - (This is the first bullet of IH when $t = \#SCCs$)

Punchline:
we can find SCCs in time $O(n + m)$

Algorithm:

- Do DFS to create a **DFS forest**.
 - Choose starting vertices in any order.
 - Keep track of finishing times.
- Reverse all the edges in the graph.
- Do DFS again to create **another DFS forest**.
 - This time, order the nodes in the reverse order of the finishing times that they had from the first DFS run.
- The SCCs are the different trees in the **second DFS forest**.



(Clearly it wasn't obvious since it took all class to do! But hopefully it is less mysterious now.)

Recap

- Depth First Search reveals a very useful structure!
 - We saw Monday that this structure can be used to do **Topological Sorting** in time $O(n+m)$
 - Today we saw that it can also find **Strongly Connected Components** in time $O(n + m)$
 - This was pretty non-trivial.

Next time

- Dijkstra's algorithm!

BEFORE Next time

- Pre-lecture exercise: weighted graphs!