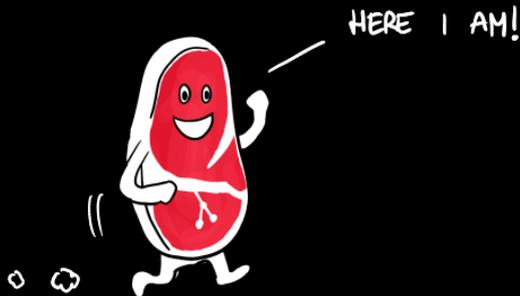


THEORETICAL ANALYSIS OF GIT BISECT

Julien COURTIEL (Université de Caen Normandie)

with Paul DORBEC and Romain LECOQ (Université de Caen Normandie)

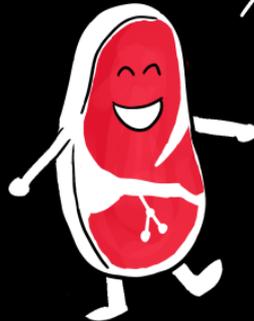


Journées APR (30 Juin 2024)

PART 1 LOOKING FOR THE FAULTY COMMIT

ARE WE SURE TO KEEP ON DOING THAT "JOKE"?

YES, I'M VERY SUBVERSIVE!



DID YOU CALL ME?



QUESTION: What do you use to share files with your coauthors?

THE TIER LIST

QUESTION: What do you use to share files with your coauthors?

god-like



good



trash



spawn of the devil



THE TIER LIST

QUESTION: What do you use to share files with your coauthors?

god-like



good



trash



spawn of the devil



GIT AND ITS COMMIT GRAPHS



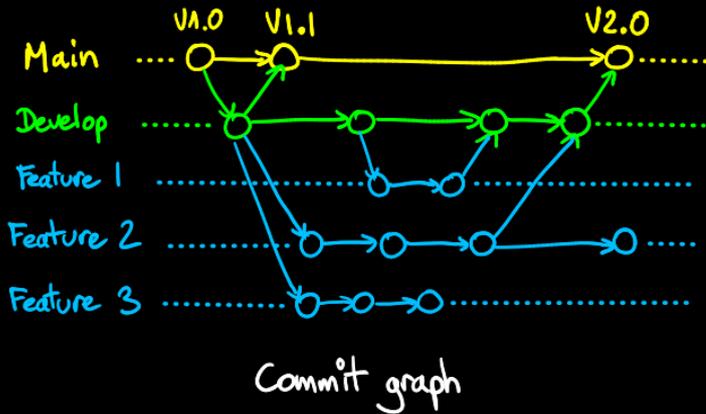
is a distributed version control system where the revisions (or "commits") are arranged as a Directed Acyclic Graph (DAG)



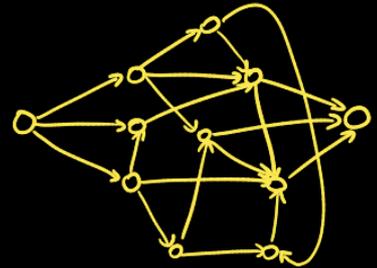
Commit graph

GIT AND ITS COMMIT GRAPHS

 is a distributed version control system where the revisions (or "commits") are arranged as a Directed Acyclic Graph (DAG)

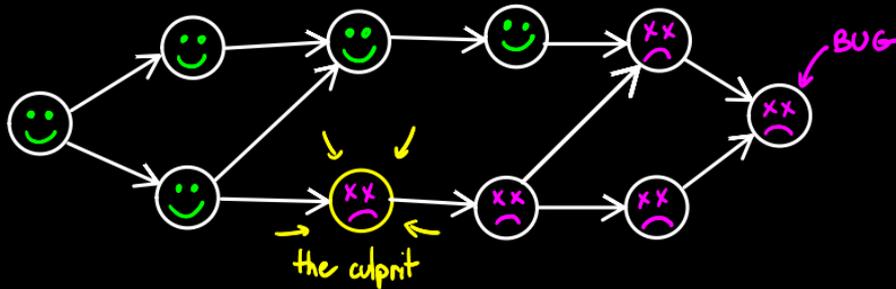


or



Also a legit
commit graph!

PROBLEM: FINDING A REGRESSION



Input

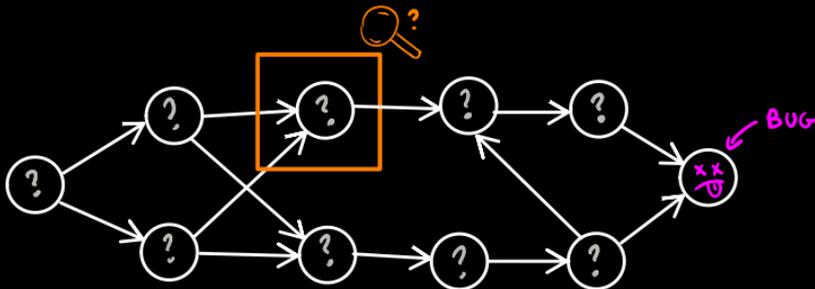
A commit graph in which a commit is known to be bugged, the other commits are bugged or clean (=bug-free)

Question

Which commit has originally introduced the bug?

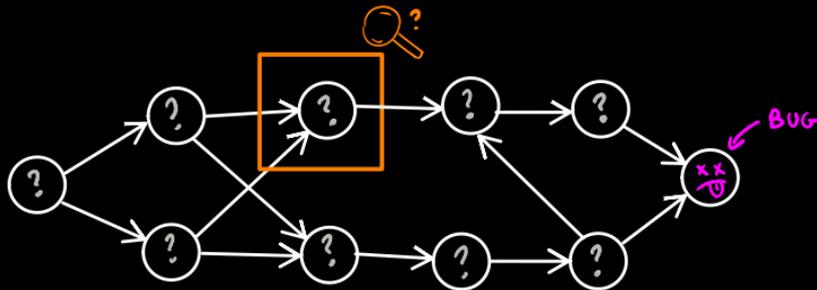
HOW TO INVESTIGATE

Unique operation: QUERY of a commit with unknown status



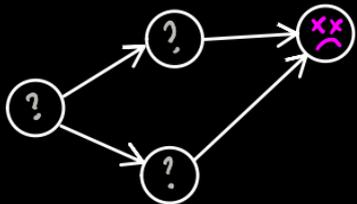
HOW TO INVESTIGATE

Unique operation: QUERY of a commit with unknown status



If bugged,

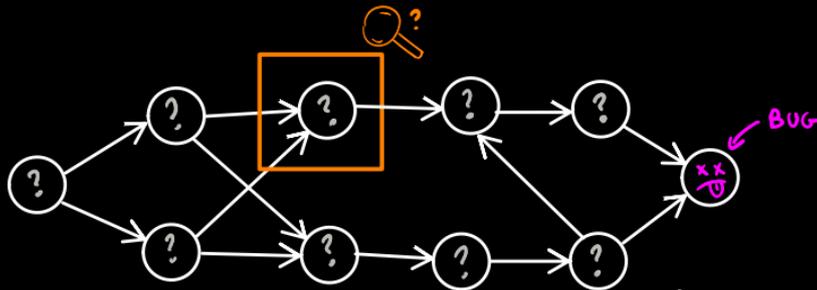
then the faulty commit is an ancestor of this commit



ancestor of a vertex v =
 v or
an ancestor of a parent of v

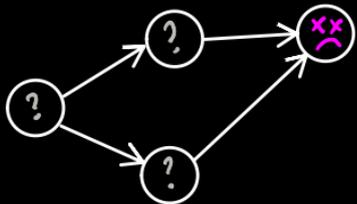
HOW TO INVESTIGATE

Unique operation: QUERY of a commit with unknown status



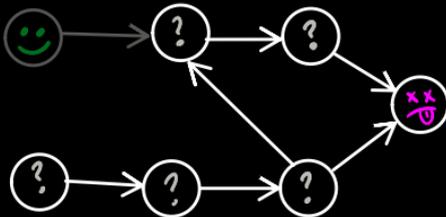
If bugged,

then the faulty commit is an ancestor of this commit



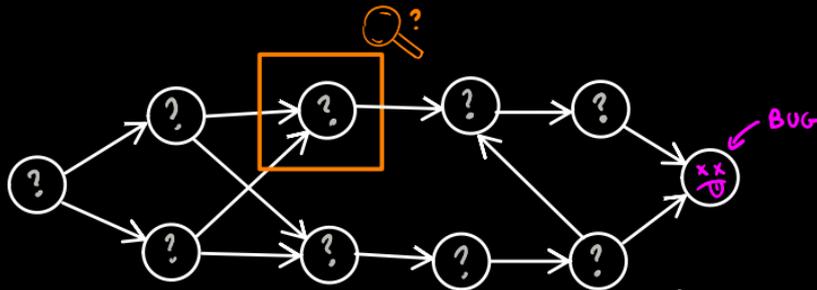
If clean,

then the faulty commit is not an ancestor of this commit



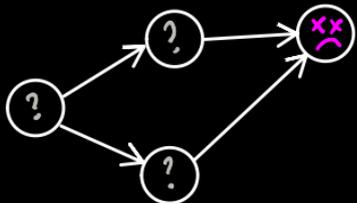
HOW TO INVESTIGATE

Unique operation: QUERY of a commit with unknown status



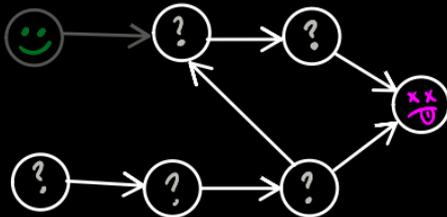
If bugged,

then the faulty commit is an ancestor of this commit



If clean,

then the faulty commit is not an ancestor of this commit



The faulty commit is found whenever there remains only 1 suspect.

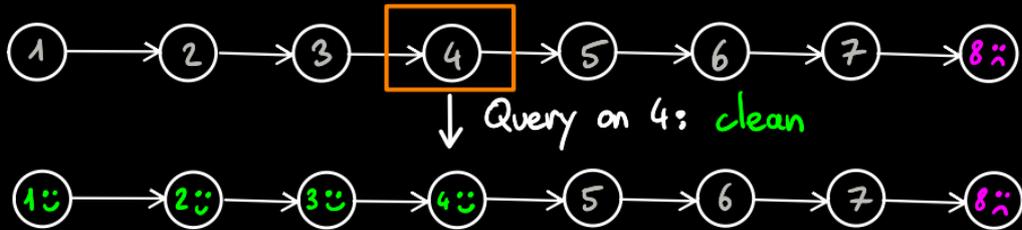
FIRST EXAMPLE : CHAINS



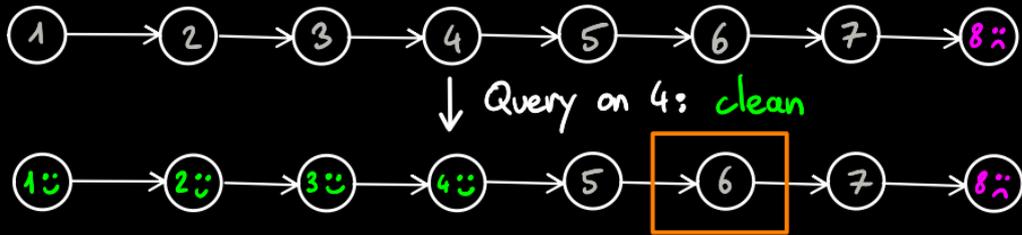
FIRST EXAMPLE : CHAINS



FIRST EXAMPLE : CHAINS



FIRST EXAMPLE : CHAINS



FIRST EXAMPLE : CHAINS



↓ Query on 4: clean



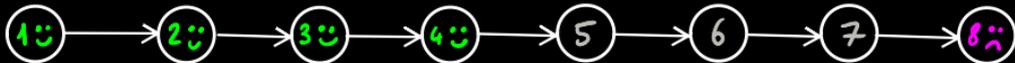
↓ Query on 6: bugged



FIRST EXAMPLE : CHAINS



↓ Query on 4: clean



↓ Query on 6: bugged



FIRST EXAMPLE : CHAINS



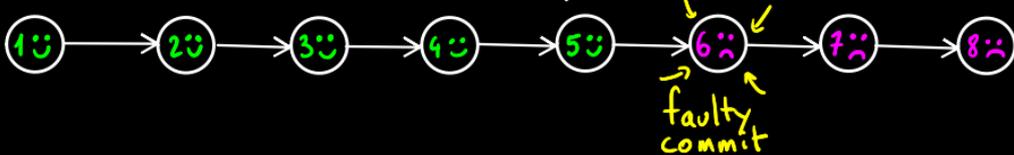
↓ Query on 4: clean



↓ Query on 6: bugged



↓ Query on 5: clean



FIRST EXAMPLE : CHAINS



↓ Query on 4: clean



↓ Query on 6: bugged



↓ Query on 5: clean

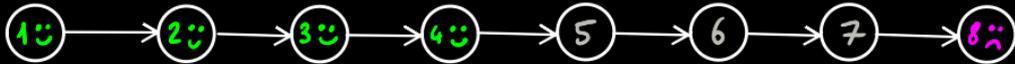


optimal strategy = binary search

FIRST EXAMPLE : CHAINS



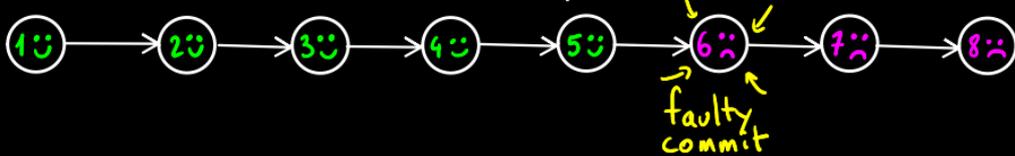
↓ Query on 4: clean



↓ Query on 6: bugged



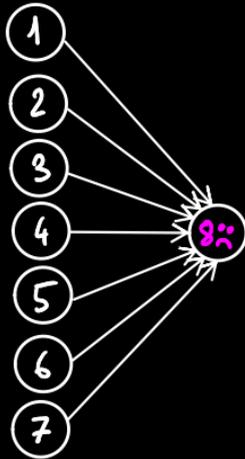
↓ Query on 5: clean



optimal strategy = binary search

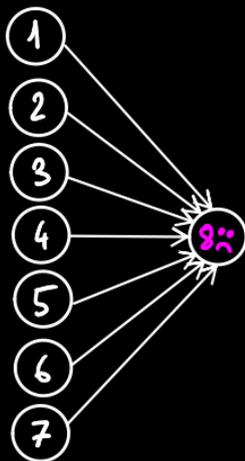
More generally, number of queries in an optimal strategy in a chain of length $n = \lceil \log_2(n) \rceil$

SECOND EXAMPLE : OCTOPUSES



optimal strategy =

SECOND EXAMPLE : OCTOPUSES

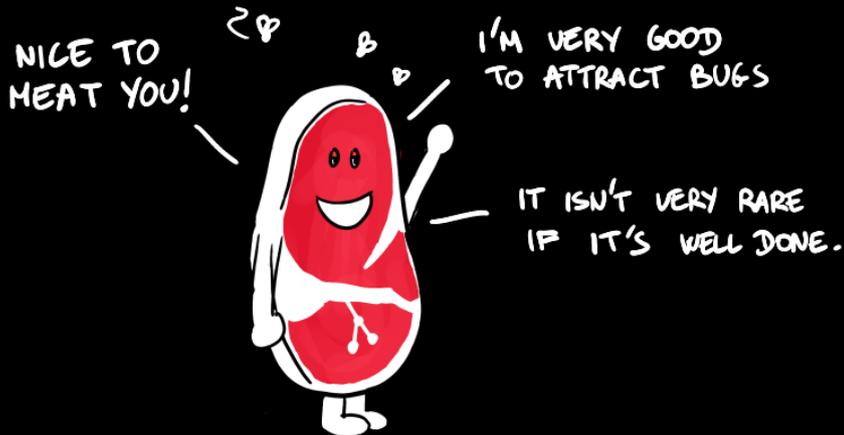


optimal strategy = whatever

More generally, number of queries in an optimal strategy in an octopus of size $n = n - 1$

PART II

GIT BISECT

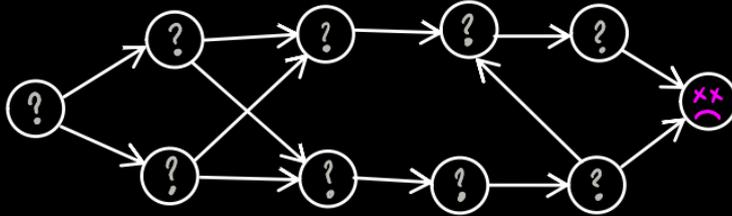


THE GIT BISECT ALGORITHM



git uses a heuristic to find the faulty commit: git bisect

STEP 1:



STEP 2 :

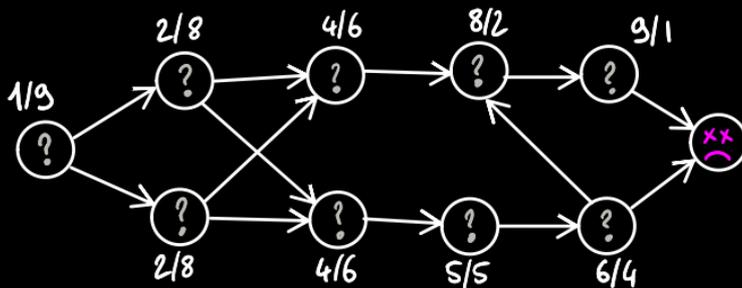
STEP 3 :

THE GIT BISECT ALGORITHM



uses a heuristic to find the **faulty commit**: git bisect

STEP 1: Compute the number of ancestors / non-ancestors for each commit



STEP 2 :

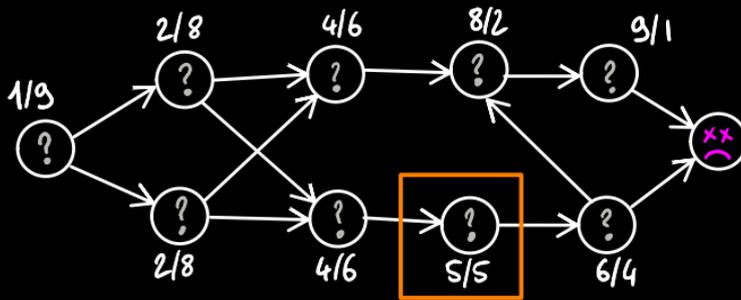
STEP 3 :

THE GIT BISECT ALGORITHM



uses a heuristic to find the **faulty commit**: git bisect

STEP 1: Compute the number of ancestors / non-ancestors for each commit



STEP 2: **Query** the vertex with the most balanced ratio

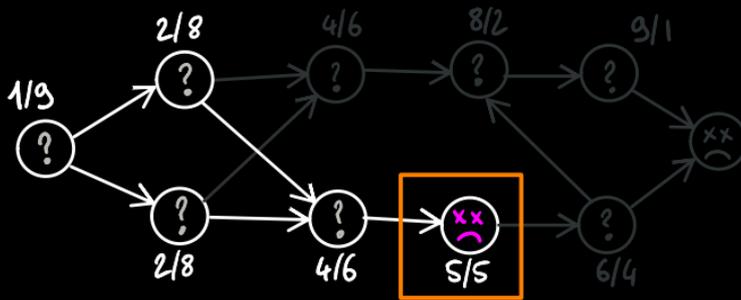
STEP 3:

THE GIT BISECT ALGORITHM



uses a heuristic to find the **faulty commit**: git bisect

STEP 1: Compute the number of ancestors / non-ancestors for each commit



STEP 2: **Query** the vertex with the most balanced ratio

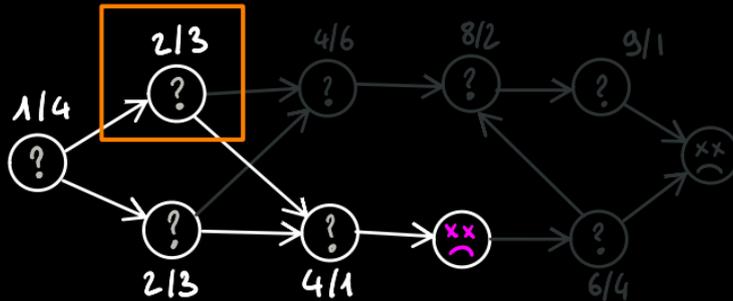
STEP 3: Delete the innocent commits and recurse.

THE GIT BISECT ALGORITHM



uses a heuristic to find the **faulty commit**: git bisect

STEP 1: Compute the number of ancestors / non-ancestors for each commit



STEP 2: Query $\frac{a}{b}$ the vertex with the most balanced ratio

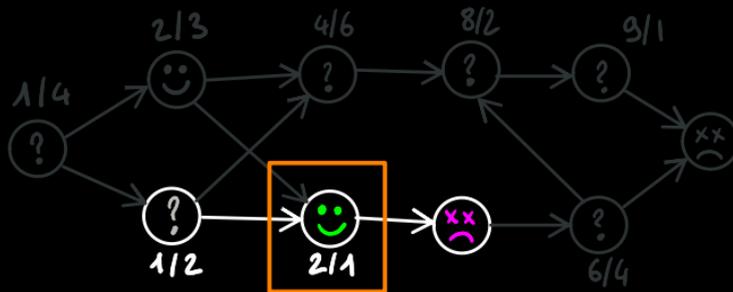
STEP 3: Delete the innocent commits and recurse.

THE GIT BISECT ALGORITHM



uses a heuristic to find the **faulty commit**: git bisect

STEP 1: Compute the number of ancestors / non-ancestors for each commit



STEP 2: Query $\frac{a}{b}$ the vertex with the most balanced ratio

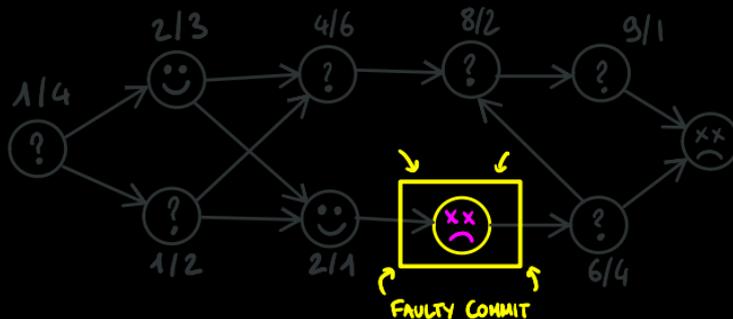
STEP 3: Delete the innocent commits and recurse.

THE GIT BISECT ALGORITHM



git uses a heuristic to find the **faulty commit**: git bisect

STEP 1: Compute the number of ancestors / non-ancestors for each commit



STEP 2: Query $\frac{a}{b}$ the vertex with the most balanced ratio

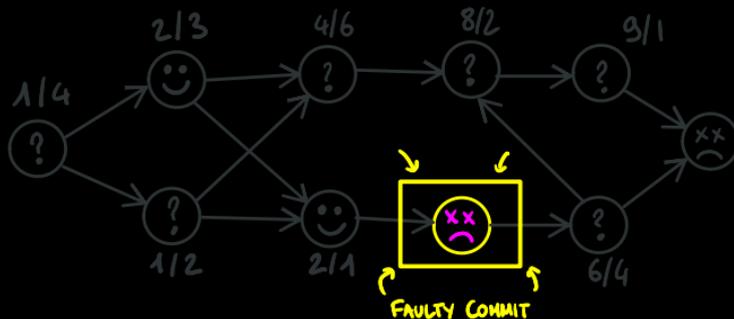
STEP 3: Delete the innocent commits and recurse.

THE GIT BISECT ALGORITHM



git uses a heuristic to find the faulty commit: git bisect

STEP 1: Compute the number of ancestors / non-ancestors for each commit



STEP 2: Query the vertex with the most balanced ratio a

STEP 3: Delete the innocent commits and recurse.

QUESTION: Does git bisect always give an optimal strategy?

HOW GOOD IS GIT BISECT?

QUESTION: Does `git bisect` always give an optimal strategy?

HOW GOOD IS GIT BISECT?

QUESTION: Does **git bisect** always give an **optimal strategy**?

NO The Regression Search Problem is NP-complete.

HOW GOOD IS GIT BISECT?

QUESTION: Does **git bisect** always give an **optimal strategy**?

NO The Regression Search Problem is NP-complete.

But is it really that bad?

HOW GOOD IS GIT BISECT?

QUESTION: Does **git bisect** always give an **optimal strategy**?

NO The Regression Search Problem is NP-complete.

But is it really that bad? **YEAH**

HOW GOOD IS GIT BISECT?

QUESTION: Does **git bisect** always give an **optimal strategy**?

NO The Regression Search Problem is NP-complete.

But is it really that bad? **YEAH!**

Proposition

For any k , there exists a DAG such that
an **optimal strategy** uses k queries
and **git bisect** always uses $2^{k-1} - 1$ queries.

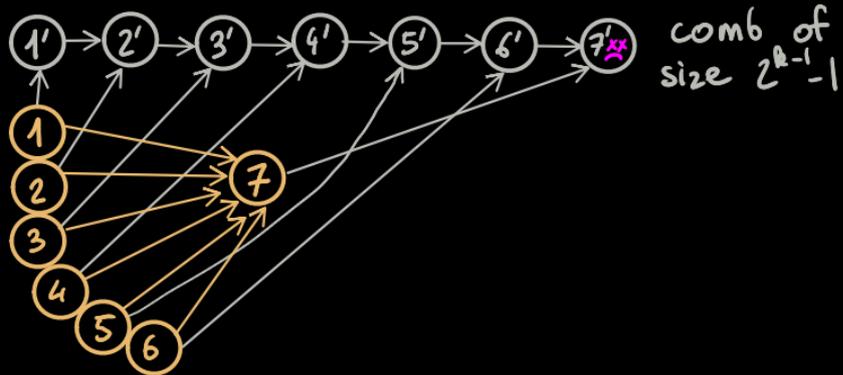
HOW GOOD IS GIT BISECT?

Proposition

For any k , there exists a DAG such that an optimal strategy uses k queries and `git bisect` always uses $2^{k-1} - 1$ queries.

Proof for $k=4$:

octopus of size $2^{k-1} - 1$

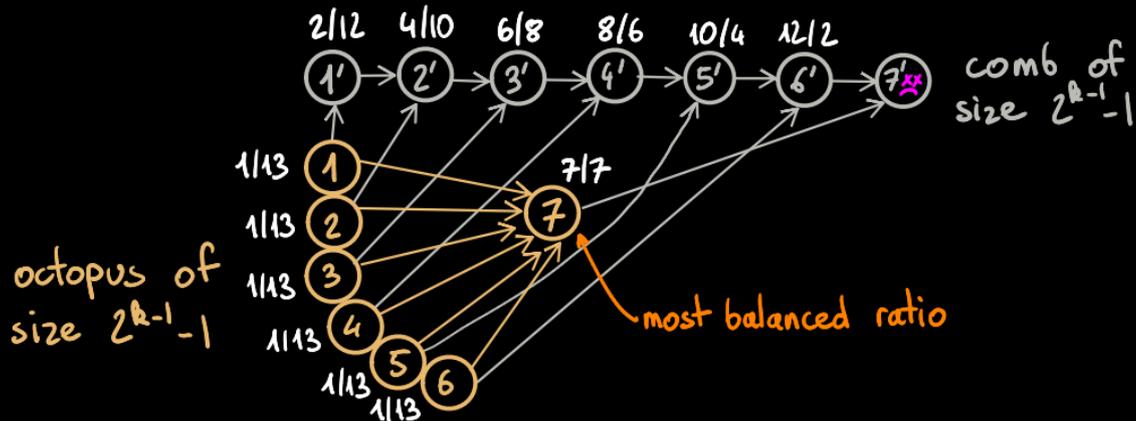


How GOOD IS GIT BISECT?

Proposition

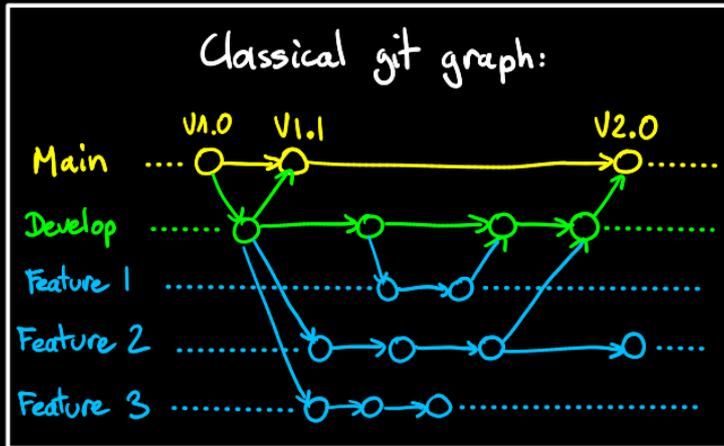
For any k , there exists a DAG such that an optimal strategy uses k queries and **git bisect** always uses $2^{k-1} - 1$ queries.

Proof for $k=4$:



BACK TO REALITY?

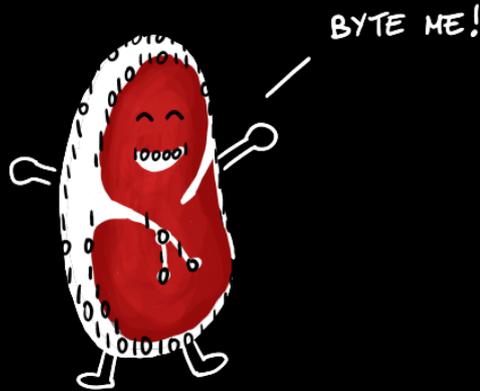
Octopus substructures are unrealistic



Usually, we never merge more than 2 branches.

(Otherwise it is called an octopus merge )

PART III - GIT BISECT ON BINARY DAGs

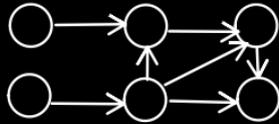


BINARY DAG

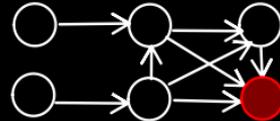
Definition

binary DAG = DAG where the vertices have indegree ≤ 2

Ex:



Good



Bud

Theorem

git bisect is a $\frac{1}{\log_2(\frac{3}{2})}$ -approximation algorithm when it is used on binary DAGs.

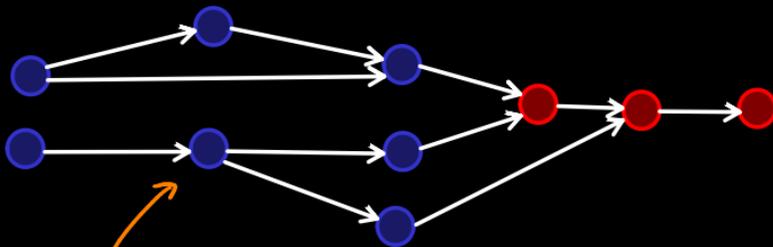
$\frac{1}{\log_2(\frac{3}{2})} \approx 1.71$ is the optimal constant.

EXISTENCE OF A BALANCED VERTEX

Lemma

Any binary DAG with n vertices has a vertex x such that

$$\frac{n}{3} \leq \text{nb of ancestors of } x \leq \frac{2n+1}{3}$$



n ? Where is it??

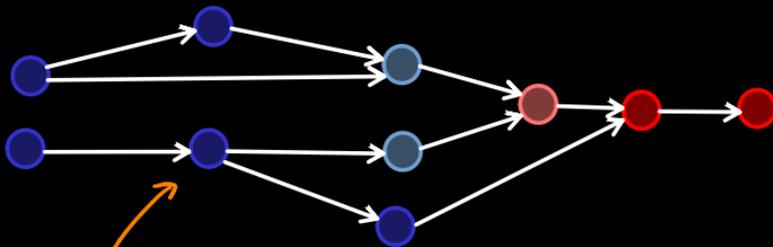
- = more ancestors than non-ancestors
- = more non-ancestors than ancestors

EXISTENCE OF A BALANCED VERTEX

Lemma

Any binary DAG with n vertices has a vertex x such that

$$\frac{n}{3} \leq \text{nb of ancestors of } x \leq \frac{2n+1}{3}$$



n. Where is it?

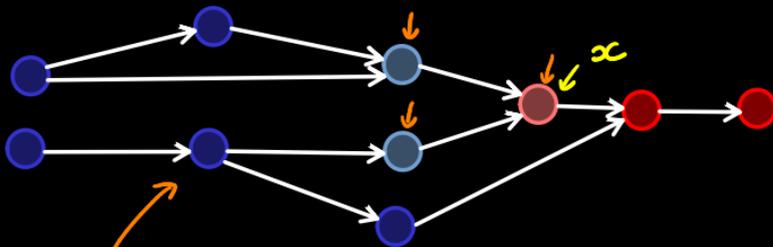
- = more ancestors than non-ancestors
- = more non-ancestors than ancestors
- = ● where all parents are ●
- = ● parent of a ●

EXISTENCE OF A BALANCED VERTEX

Lemma

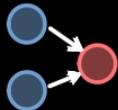
Any binary DAG with n vertices has a vertex x such that

$$\frac{n}{3} \leq \text{nb of ancestors of } x \leq \frac{2n+1}{3}$$



n. Where is it?

It's in any



- = more ancestors than non-ancestors
- = more non-ancestors than ancestors
- = ● where all parents are ●
- = ● parent of a ●

EXISTENCE OF A BALANCED VERTEX

Lemma

Any binary DAG with n vertices has a vertex x such that

$$\frac{n}{3} \leq \text{nb of ancestors of } x \leq \frac{2n+1}{3}$$

Proof of the $\frac{1}{\log_2(\frac{3}{2})}$ - approximation:

At each query, **git bisect** eliminates $\geq 1/3$ commits by choosing x or a better vertex



- number of **git bisect** queries $\approx \log_{\frac{3}{2}}(n)$
- **optimal** number of queries $\geq \log_2(n)$

PROBLEMATIC BINARY DAGs FOR GIT BISECT

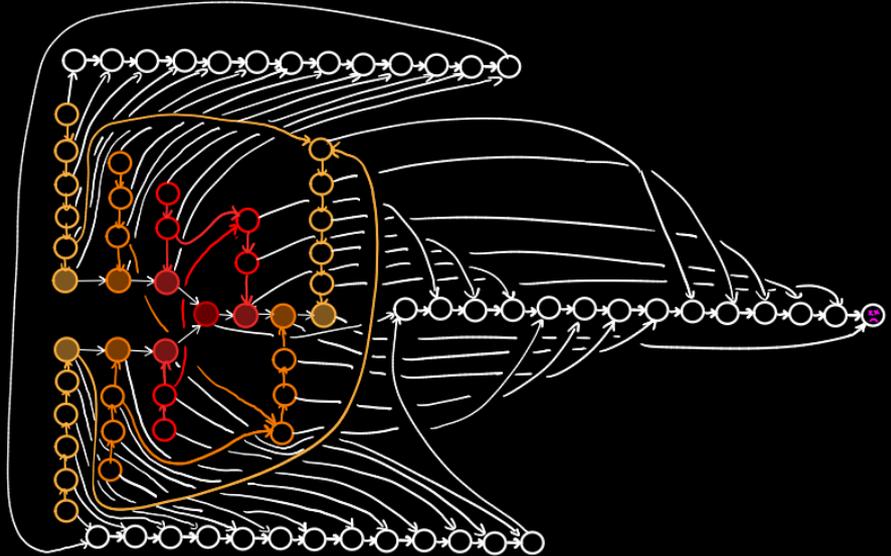
For every $k > 0$, there exists a binary DAG with:

$\approx 6k \left(\frac{3}{2}\right)^k$ vertices,

$\geq k$ **git bisect** queries,

$\approx \log_2 \left(\frac{3}{2}\right) \times k$ queries in the **optimal** strategy.

e.g. for $k=3$:



PROBLEMATIC BINARY DAGs FOR GIT BISECT

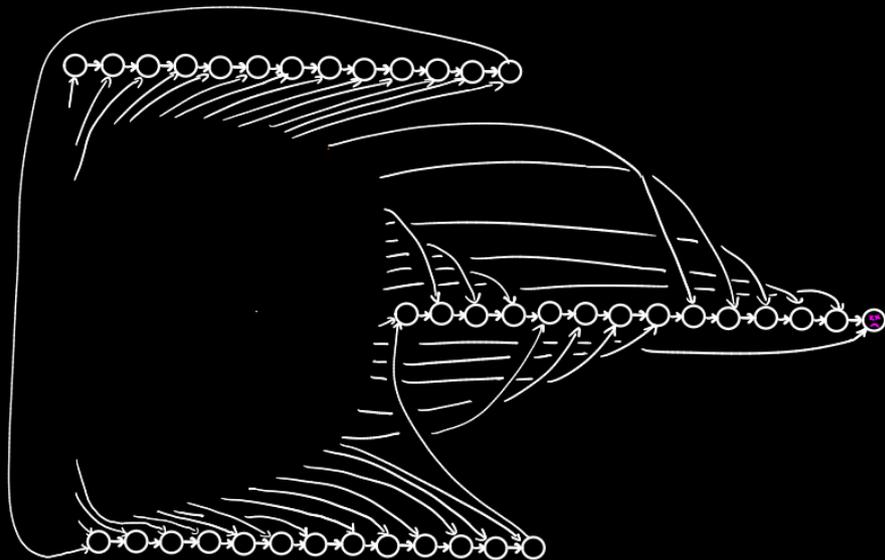
For every $k > 0$, there exists a binary DAG with:

$\approx 6k \left(\frac{3}{2}\right)^k$ vertices,

$\geq k$ **git bisect** queries,

$\approx \log_2\left(\frac{3}{2}\right) \times k$ queries in the **optimal** strategy.

e.g. for $k=3$:



PROBLEMATIC BINARY DAGs FOR GIT BISECT

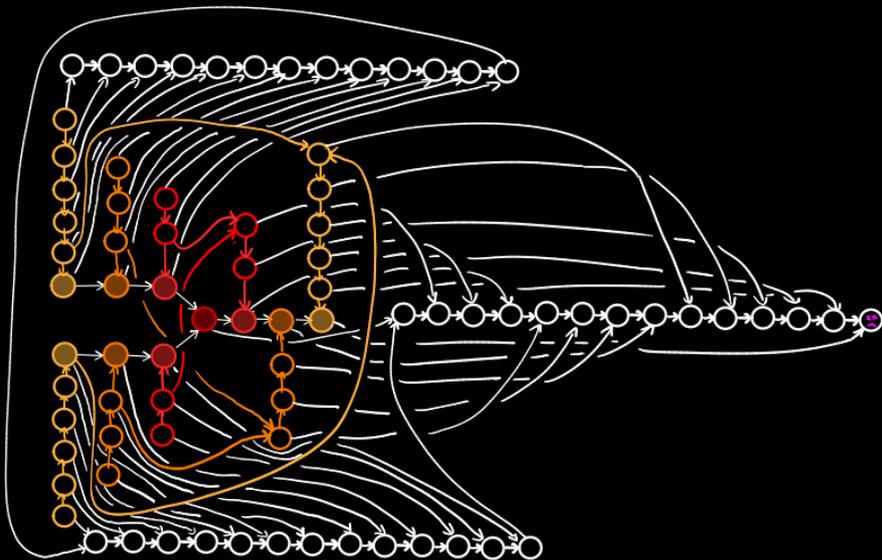
For every $k > 0$, there exists a binary DAG with:

$\approx 6k \left(\frac{3}{2}\right)^k$ vertices,

$\geq k$ **git bisect** queries,

$\approx \log_2 \left(\frac{3}{2}\right) \times k$ queries in the **optimal** strategy.

e.g. for $k=3$:



PROBLEMATIC BINARY DAGs FOR GIT BISECT

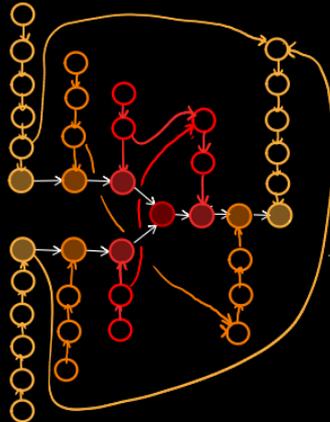
For every $k > 0$, there exists a binary DAG with:

$\approx 6k \left(\frac{3}{2}\right)^k$ vertices,

$\geq k$ **git bisect** queries,

$\approx \log_2 \left(\frac{3}{2}\right) \times k$ queries in the **optimal** strategy.

e.g. for $k=3$:



A BETTER ALGORITHM? GOLDEN BISECT

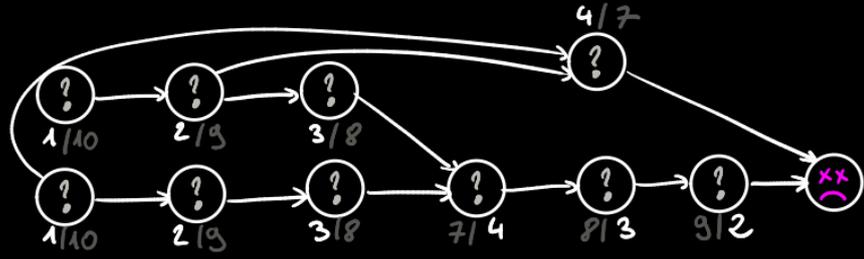
theoretical analysis of `git bisect` → new (better?) algorithm

A BETTER ALGORITHM? GOLDEN BISECT

theoretical analysis of **git bisect** → new (better?) algorithm

★ GOLDEN BISECT ★

Step 1: Compute the number of ancestors/number of non-ancestors for each vertex and keep the smaller of the 2 numbers.

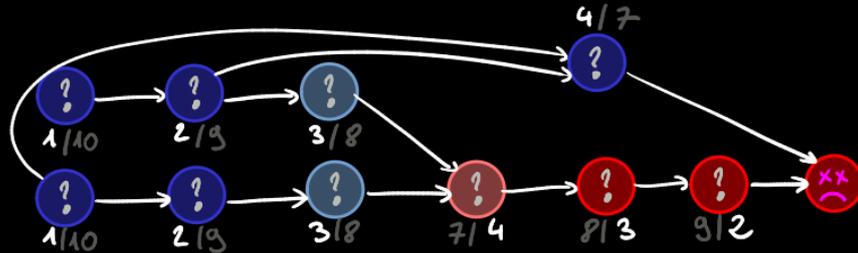


A BETTER ALGORITHM? GOLDEN BISECT

theoretical analysis of **git bisect** → new (better?) algorithm

★ GOLDEN BISECT ★

Step 1: Compute the number of ancestors/number of non-ancestors for each vertex and keep the smaller of the 2 numbers.



Step 2: If the largest number is $\geq \frac{\text{nb vertices}}{\phi^2} \approx 0.38 \times \text{nb vertices}$,
query a commit with this number.

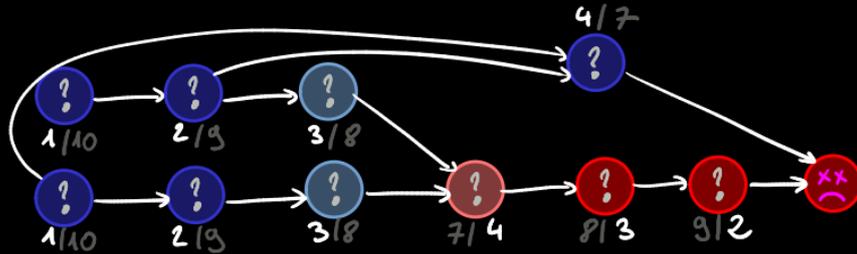
$$\phi = \text{golden ratio} = \frac{1 + \sqrt{5}}{2}$$

A BETTER ALGORITHM? GOLDEN BISECT

theoretical analysis of **git bisect** → new (better?) algorithm

★ GOLDEN BISECT ★

Step 1: Compute the number of ancestors/number of non-ancestors for each vertex and keep the smaller of the 2 numbers.



Step 2: If the largest number is $\geq \frac{\text{nb vertices}}{\phi^2} \approx 0.38 \times \text{nb vertices}$,
query a commit with this number.

$$\phi = \text{golden ratio} = \frac{1 + \sqrt{5}}{2}$$

Otherwise, query a ● or a ●
with the largest number

- = more ancestors than non-ancestors
- = more non-ancestors than ancestors
- = ● where all parents are ●
- = ● parent of a ●

Step 3: Delete the innocent commits and recurse.

THEORETICAL ANALYSIS OF GOLDEN BISECT

Theorem

golden bisect is a $\frac{1}{\log_2(\phi)}$ - approximation algorithm for binary DAGs, where $\phi = \text{golden ratio}$.

$\frac{1}{\log_2(\phi)} \approx 1,44$ is the optimal constant.

But sometimes **git bisect** is better than golden bisect!

OPEN QUESTIONS

→ Is the Regression Search Problem NP-complete when restricted to binary DAGs?

(We proved that a variant, the Confined Regression Search Problem is NP-complete for binary DAGs)

→ study of **git bisect** in average case?

but what is a random "git graph"?

(\Rightarrow definition of git graph in [Lecoq's PhD],
random generation in [C., Pépin])

→ Is **git-bisect** a 2-approximation algorithm for trees?

THANK YOU!



STEAK UN
AU REVOIR!