# Theoretical Analysis of Git Bisect
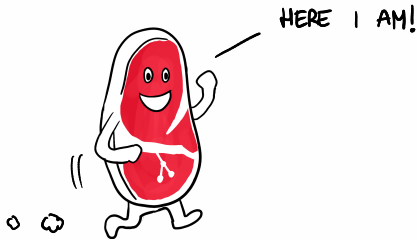
Julien COURTIEL (Université de Caen Normandie)
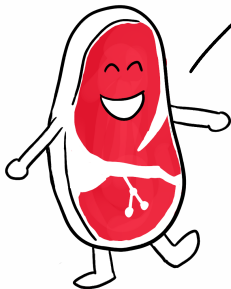
with Paul DORBEC and Romain LECOQ (Université de Caen Normandie)

HERE I AM!

I SAID "GIT BISECT", NOT "GIT BEEFSTEAK"!

SIMPLE QUESTION: WHAT DO YOU USE TO SHARE FILES WITH YOUR COAUTHORS?
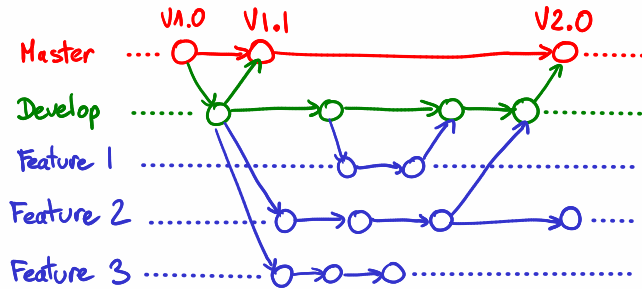
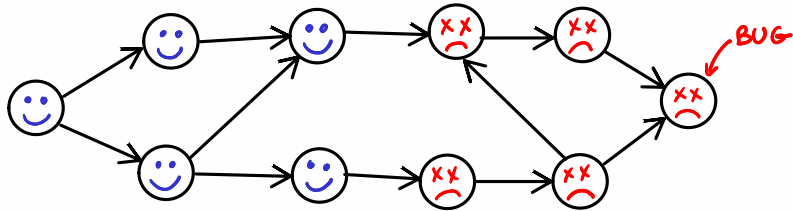Different ways to share files within a project

# GIT AND ITS COMMIT GRAPH

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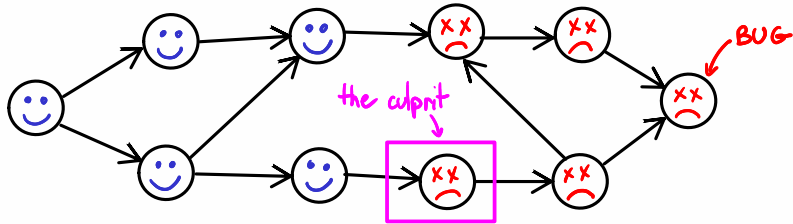


Commit graph

# PROBLEM : FINDING THE SOURCE OF A BUG



**Input** A commit graph in which a commit is known to be bugged, the other commits may be bugged or bug-free.

**Question** Which commit has originally introduced the bug?

# PROBLEM : FINDING THE SOURCE OF A BUG



Input: A commit graph in which a commit is known to be bugged, the other commits may be bugged or bug-free.

Question: Which commit has originally introduced the bug?

# PROBLEM : FINDING THE SOURCE OF A BUG



**Input** A commit graph in which a commit is known to be bugged, the other commits may be bugged or bug-free.

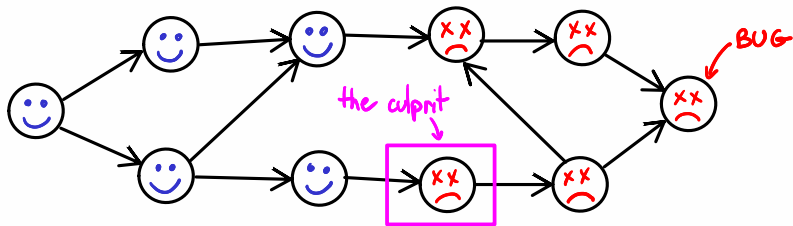**Question** Which commit has originally introduced the bug?

# PROBLEM : FINDING THE SOURCE OF A BUG



**Input** A commit graph in which a commit is known to be bugged, the other commits may be bugged or bug-free.

**Question** Which commit has originally introduced the bug?

**Assumptions**
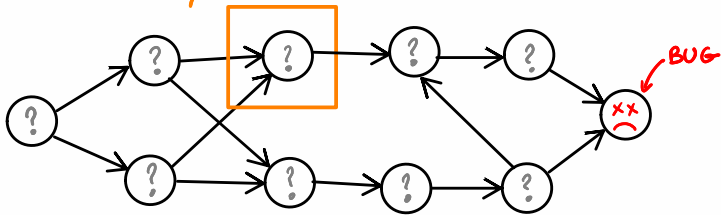- If a parent of a commit is bugged, then the commit is bugged. (Monotonous hypothesis)
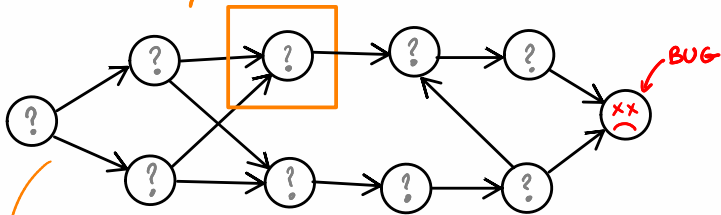- Only one commit has introduced the bug, namely the original bug.

# HOW TO CATCH THE FIRST BUG

Only operation : Query of a commit with unknown status
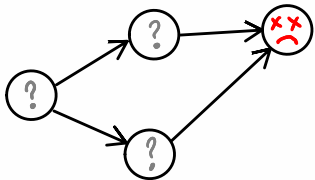
# HOW TO CATCH THE FIRST BUG

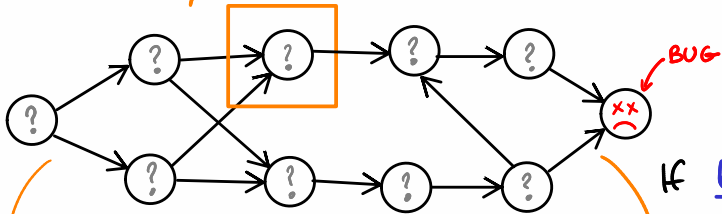Only operation : Query of a commit with unknown status



If bugged,

then the original bug is an ancestor of this commit

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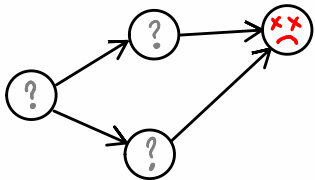$v$ or
an ancestor of a parent of $v$

# HOW TO CATCH THE FIRST BUG

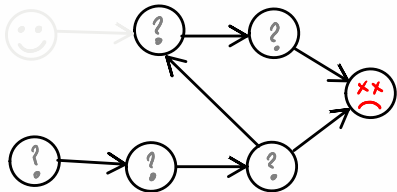Only operation : Query of a commit with unknown status



If bugged,

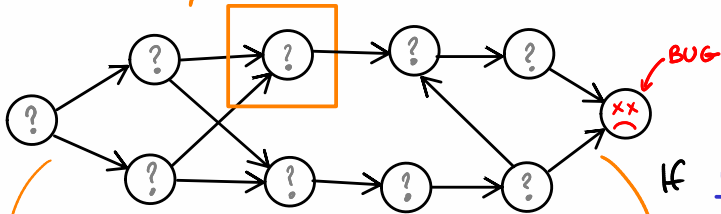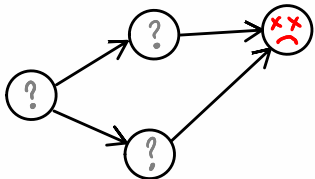then the original bug is an ancestor of this commit

If bug-free,

then the original bug is not an ancestor of this commit

# HOW TO CATCH THE FIRST BUG

Only operation : Query of a commit with unknown status



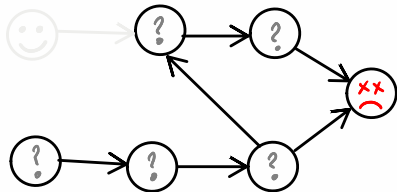If **bugged**,
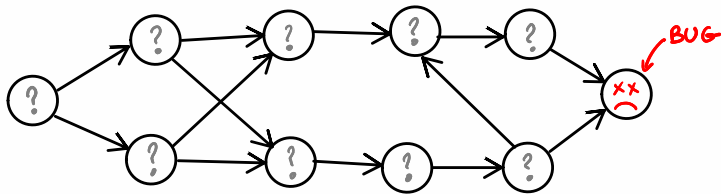
then the **original bug** is an ancestor of this commit

If **bug-free**,

then the **original bug** is **not** an ancestor of this commit

The **original bug** is found whenever the remaining graph has only 1 vertex

# PRECISE DEFINITION OF THE PROBLEM

<u>Input</u>: a DAG where each vertex has an <span style="color:gray">unknown</span> status, except one, which is <span style="color:red">bugged</span>, such that every vertex is an ancestor of this <span style="color:red">bugged</span> vertex



<u>Output</u>: A <span style="color:purple">strategy</span> that finds the <span style="color:magenta">original bug</span> with a minimal number of <span style="color:orange">queries</span> in the worst-case scenario = optimal strategy
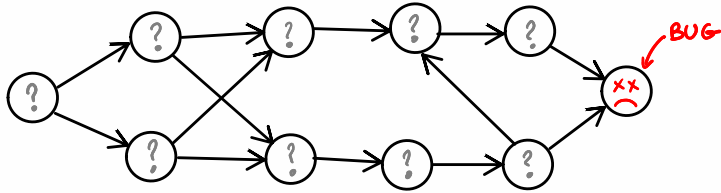
# PRECISE DEFINITION OF THE PROBLEM
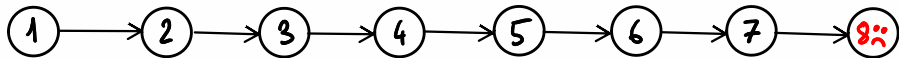
<u>Input</u>: a DAG where each vertex has an unknown status, except one, which is bugged, such that every vertex is an ancestor of this bugged vertex
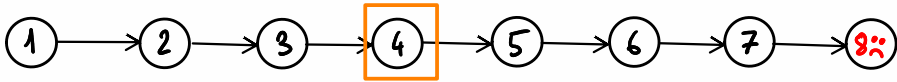


<u>Output</u>: A strategy that finds the original bug with a minimal number of queries in the worst-case scenario = optimal strategy
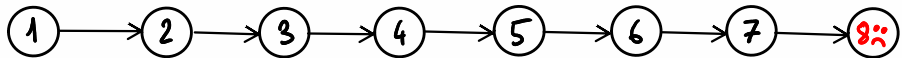
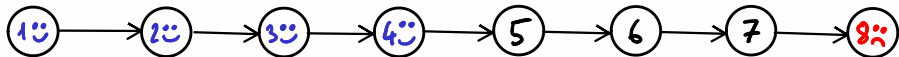In real life, queries are costly.

# FIRST EXAMPLE: A CHAIN

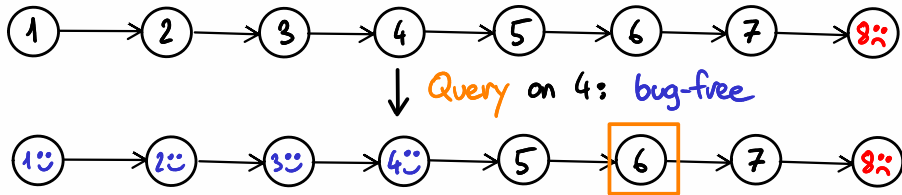$1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5 \rightarrow 6 \rightarrow 7 \rightarrow$ 😞

# FIRST EXAMPLE: A CHAIN
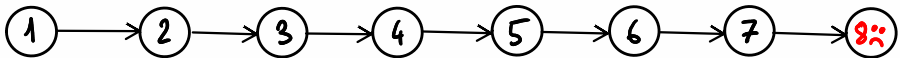
# FIRST EXAMPLE: A CHAIN

# FIRST EXAMPLE:   A CHAIN



① → ② → ③ → ④ → ⑤ → ⑥ → ⑦ → ⑧

↓ Query on 4: bug-free

①☺ → ②☺ → ③☺ → ④☺ → ⑤ → ⑥ → ⑦ → ⑧

# FIRST EXAMPLE:  A CHAIN



① → ② → ③ → ④ → ⑤ → ⑥ → ⑦ → ⑧

Query on 4: bug-free

①☺ → ②☺ → ③☺ → ④☺ → ⑤ → ⑥ → ⑦ → ⑧☹

Query on 6: bugged

①☺ → ②☺ → ③☺ → ④☺ → ⑤ → ⑥☹ → ⑦☹ → ⑧☹

# FIRST EXAMPLE: A CHAIN

# FIRST EXAMPLE: A CHAIN



**Optimal strategy** = binary search

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

# SECOND EXAMPLE: A RAKE

# SECOND EXAMPLE: A RAKE



$\underline{\text{Optimal strategy}}$ = whatever

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

# STRATEGY TREE



Strategy tree for binary search:

# STRATEGY TREE

Strategy tree for binary search:



In short:

# STRATEGY TREE



Strategy tree for binary search:

height of a strategy tree = number of requests in the worst-case scenario

# COMPLEXITY OF THE PROBLEM

Finding the number of queries in an optimal strategy is ...

- <u>NP-complete</u> for general DAGs [Carmo Donadelli Kohayakawa Laber 2004]

    Certificate: Strategy tree
    Reduction to: Cover by 3-sets

- but <u>polynomial</u> for trees ...
    [Ben-Asher Farchi Newman 2000]



... more precisely, linear. [Mozes Onak Weimann 2008]

# LET ME INTRODUCE YOU GIT BISECT

git uses a heuristic algorithm to find the original bug: git bisect

originally written by Linus Torvalds himself

now maintained by Junio Hamano

it's him

In the source code of git bisect:

```
/*
 * This is a truly stupid algorithm, but it's only
 * used for bisection, and we just don't care enough.
 *
```

# LET ME INTRODUCE YOU GIT BISECT

git uses a heuristic algorithm to find the original bug: git bisect

<u>Step 1</u>: Compute the number of ancestors/number of non-ancestors for each vertex:

# LET ME INTRODUCE YOU GIT BISECT

**git** uses a heuristic algorithm to find the *original bug*: *git bisect*

<u>Step 1</u>: Compute the number of ancestors/number of non-ancestors for each vertex:



<u>Step 2</u>: **Query** on a vertex with the most balanced ratio
(max of both numbers)

<u>Step 3</u>: Recurse

# LET ME INTRODUCE YOU GIT BISECT

git uses a heuristic algorithm to find the original bug: git bisect

<u>Step 1</u>: Compute the number of ancestors/number of non-ancestors for each vertex:



<u>Step 2</u>: Query on a vertex with the most balanced ratio
(max of both numbers)

<u>Step 3:</u> Recurse

# LET ME INTRODUCE YOU GIT BISECT

**git** uses a heuristic algorithm to find the *original bug*: git bisect

<u>Step 1</u>: Compute the number of ancestors/number of non-ancestors for each vertex:



<u>Step 2</u>: Query on a vertex with the most balanced ratio
(max of both numbers)

<u>Step 3:</u> Recurse

# LET ME INTRODUCE YOU GIT BISECT

git uses a heuristic algorithm to find the *original bug*: git bisect

<u>Step 1</u>: Compute the number of ancestors/number of non-ancestors for each vertex:



<u>Step 2</u>: Query on a vertex with the most balanced ratio
(max of both numbers)

<u>Step 3:</u> Recurse

# LET ME INTRODUCE YOU GIT BISECT

**git** uses a heuristic algorithm to find the *original bug*: git bisect

<u>Step 1</u>: Compute the number of ancestors/number of non-ancestors for each vertex:



<u>Step 2</u>: Query on a vertex with the most balanced ratio
(max of both numbers)

<u>Step 3:</u> Recurse

# LET ME INTRODUCE YOU GIT BISECT

git uses a heuristic algorithm to find the *original bug*: git bisect

__Step 1__: Compute the number of ancestors/number of non-ancestors for each vertex:



__Step 2__: Query on a vertex with the most balanced ratio
(max of both numbers)

__Step 3__: Recurse

# LET ME INTRODUCE YOU GIT BISECT

git uses a heuristic algorithm to find the *original bug*: git bisect

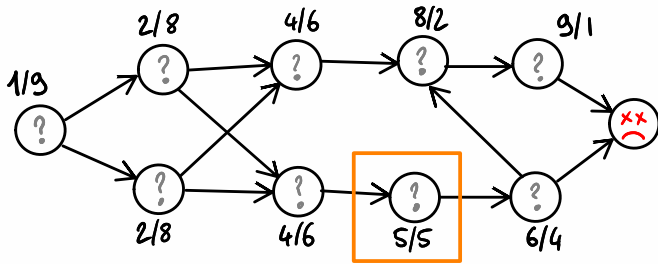<u>Step 1</u>: Compute the number of ancestors/number of non-ancestors for each vertex:



<u>Step 2</u>: Query on a vertex with the most balanced ratio
(max of both numbers)

<u>Step 3</u>: Recurse

# LET ME INTRODUCE YOU GIT BISECT

git uses a heuristic algorithm to find the **original bug**: git bisect

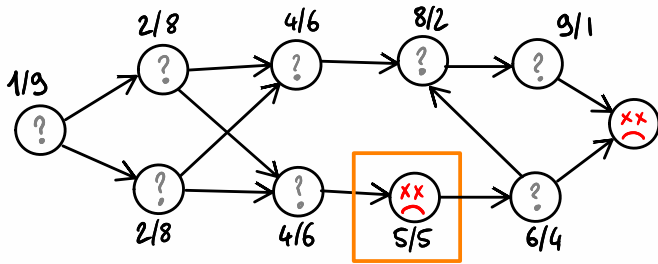<u>Step 1</u>: Compute the number of ancestors/number of non-ancestors for each vertex:



original bug

<u>Step 2</u>: Query on a vertex with the most balanced ratio
(max of both numbers)

<u>Step 3:</u> Recurse

# LET ME INTRODUCE YOU GIT BISECT

git uses a heuristic algorithm to find the original bug: git bisect

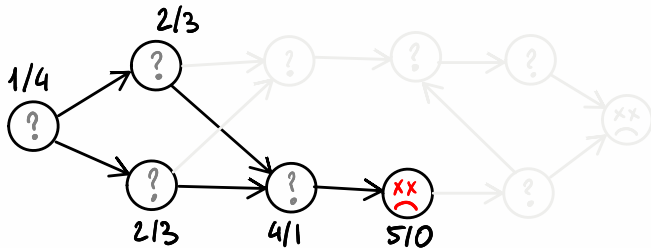Step 1: Compute the number of ancestors/number of non-ancestors for each vertex:



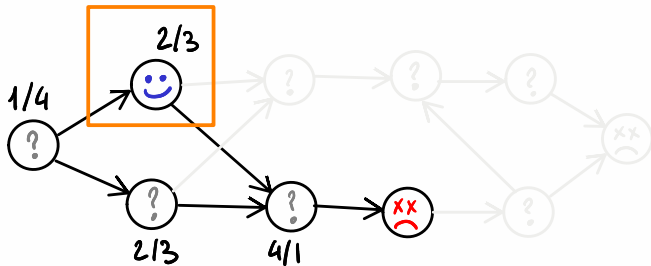original bug

Step 2: Query on a vertex with the most balanced ratio
(max of both numbers)

Step 3: Recurse

# HOW GOOD IS GIT BISECT?



Strategy tree of git bisect for this example:

Number of queries in the worst-case scenario of git bisect: 4
Number of queries of an optimal strategy: 4

# HOW GOOD IS GIT BISECT?



Strategy tree of git bisect for this example:

Number of queries in the worst-case scenario of git bisect: 4

Number of queries of an optimal strategy: 4

Question: Can git bisect always find an optimal strategy?

# HOW GOOD IS GIT BISECT?

**Question:** Can git bisect always find an optimal strategy?

**Quick answer:** absolutely no, it can't

# HOW GOOD IS GIT BISECT?

Question: Can git bisect always find an optimal strategy?

Quick answer: absolutely no, it can't

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.

# A COUNTER-EXAMPLE

| 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 $k = 4$



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

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

# A COUNTER-EXAMPLE

**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 $k = 4$



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

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

# BACK TO REALITY?

Rake substructures are unrealistic



Classical git graph:

Usually, we never merge more than 2 branches.

( Otherwise, it is called an octopus merge 🐙 )

# Part III - Git Bisect on Binary DAGs



BYTE ME!

# BINARY DAG

**Definition**

binary DAG = DAG where the vertices have indegree $\leq 2$

Ex:



Good          Bad

**Theorem**

git bisect is a $\dfrac{1}{\log_2\left(\frac{3}{2}\right)}$ - approximation algorithm

when it is used on binary DAGs.

$\dfrac{1}{\log_2\left(\frac{3}{2}\right)} \simeq 1{,}71$ is the optimal constant.

# EXISTENCE OF A BALANCED VERTEX

**Lemma**   In any binary graph of length $n$, there exists a vertex such that its number $x$ of ancestors satisfies

$$\frac{n-1}{3} < x \leq \frac{2n+1}{3}$$

# EXISTENCE OF A BALANCED VERTEX

**Lemma**    In any binary graph of length $n$, there exists a vertex such that its number $x$ of ancestors satisfies

$$\frac{n-1}{3} < x \leq \frac{2n+1}{3}$$



**Where is it?**  Consider $v =$ vertex with the least number $x$ of ancestors among those that has more ancestors than non-ancestors. ( $\frac{n}{2} \leq x$ )

The wanted vertex must be $v$ or one of its parents.

# EXISTENCE OF A BALANCED VERTEX

**Lemma** In any binary graph of length $n$, there exists a vertex such that its number $x$ of ancestors satisfies

$$\frac{n-1}{3} < x \leq \frac{2n+1}{3}$$

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

# EXISTENCE OF A BALANCED VERTEX

**Lemma** In any binary graph of length $n$, there exists a vertex such that its number $x$ of ancestors satisfies

$$\frac{n-1}{3} < x \leq \frac{2n+1}{3}$$

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

At each step, git bisect chooses a query which eliminates at least 1/3 of the vertices.

# EXISTENCE OF A BALANCED VERTEX

**Lemma**   In any binary graph of length $n$, there exists a vertex such that its number $x$ of ancestors satisfies

$$\frac{n-1}{3} < x \leq \frac{2n+1}{3}$$

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

At each step, git bisect chooses a query which eliminates at least 1/3 of the vertices.



$$\text{number of git bisect queries} \approx \log_{\frac{3}{2}}(n)$$

$$\text{optimal number of queries} \geqslant \log_2(n)$$

# TIGHTENING THE BOUND

The hard part: proving that $\frac{1}{\log_2\left(\frac{3}{2}\right)}$ is optimal

Existence of a problematic binary DAG for git bisect ?

---

**Proposition**    Let $k$ be any number.

There exists a binary DAG such that

- number of git bisect queries $= k + \lceil \log_2(k) \rceil + 2$

- optimal number of queries $\leq k \log_2\left(\frac{3}{2}\right) + \log_2(3k+6) + 4$

# STEP 1: MAXIMIZING THE NUMBER OF GB QUERIES

Binary DAG such that *git bisect* eliminates $\frac{1}{3}$ of its vertices for the $k$ first steps?

<u>Construction of an example for $k = 3$:</u>

# STEP 1: MAXIMIZING THE NUMBER OF GB QUERIES

Binary DAG such that *git bisect* eliminates $\frac{1}{3}$ of its vertices for the $k$ first steps?

Construction of an example for $k = 3$:



Step 0

$m_0 = 10$

Step 1

$m_1 = 16$

Step 2

$m_2 = 25$

Step 3

Rule    $m_i$ = nb of vertices at step $i$

Add 3 chains of length $\begin{cases} \frac{m_i + 2}{6} & \text{if } m_i \text{ is even} \\ \frac{m_i + 5}{6} & \text{if } m_i \text{ is odd} \end{cases}$

# STEP 1 : MAXIMIZING THE NUMBER OF GB QUERIES

Binary DAG such that *git bisect* eliminates $\frac{1}{3}$ of its vertices for the $k$ first steps?

Construction of an example for $k = 3$ :

Step 0

$m_0 = 10$

Step 1

$m_1 = 16$

Step 2

$m_2 = 25$

Step 3

$\simeq 3k \times \left(\frac{3}{2}\right)^k$ vertices

number of *git bisect* queries

$\geqslant k$

Rule    $m_i = $ nb of vertices at step $i$

Add 3 chains of length $\begin{cases} \dfrac{m_i + 2}{6} & \text{if } m_i \text{ is even} \\[2mm] \dfrac{m_i + 5}{6} & \text{if } m_i \text{ is odd} \end{cases}$

# STEP 1: MAXIMIZING THE NUMBER OF GB QUERIES

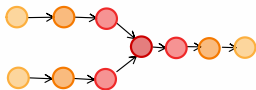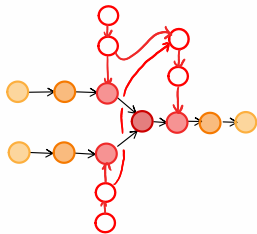Binary DAG such that git bisect eliminates $\frac{1}{3}$ of its vertices for the k first steps?
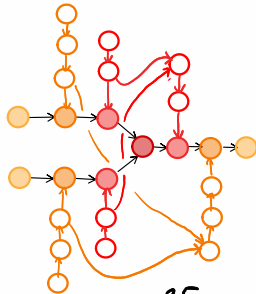
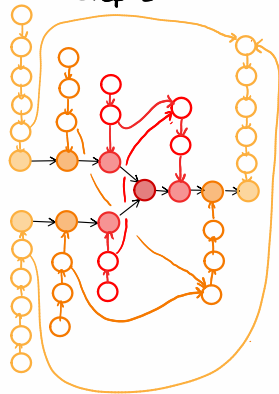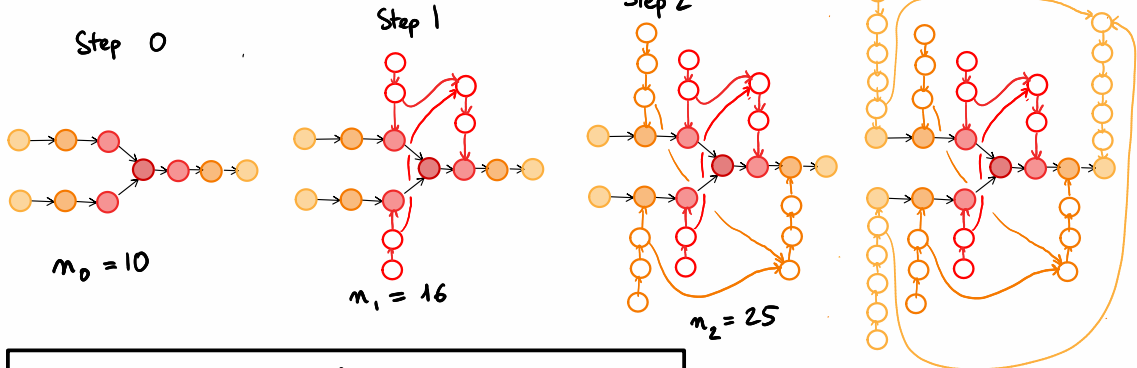Construction of an example for $k = 3$:



Step 0

gb query

Step 1

5/11    6/10

11/5

5/11

gb query

Step 2

8/17

17/8

8/17    9/16

gb query

Step 3

13/27    15/25

27/13

13/27

↑
$\simeq 3k \times \left(\frac{3}{2}\right)^k$ vertices
number of git bisect queries
$\geqslant k$

# STEP 2: TRAPPING GIT BISECT

**Theorem**

Let $D$ be a DAG with $n$ vertices, such that *git bisect* uses $gb$ queries.

There exists a DAG Comb($D$) with $2n$ vertices such that an *optimal strategy* uses $\lceil \log_2(n) \rceil + 1$ queries and if $n$ is odd, *git bisect* uses $gb+1$ queries.

**Proof:**



$D$

Comb($D$)

2/12   4/10   ....   12/2   14/0

1'  2'  3'  4'  5'  6'  7'✗✗

7/7

We assume that the vertices are topologically sorted: $1 < 2 < 3 < \ldots < 7$

# A BETTER ALGORITHM? THE GOLDEN BISECTION

theoritical analysis of git bisect → new (better?) algorithm

# A BETTER ALGORITHM?  THE GOLDEN BISECTION

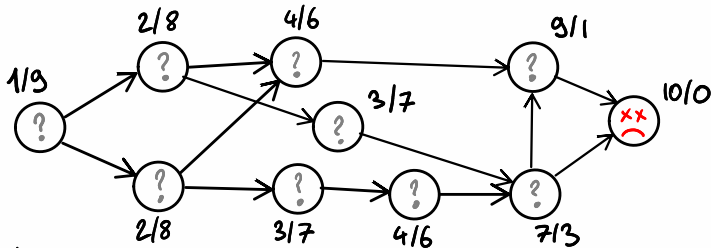theoritical analysis of git bisect → new (better?) algorithm

─────── ⭐ THE GOLDEN BISECTION ⭐ ───────

<u>Step 1</u>: Compute the number of ancestors/number of non-ancestors for each vertex:



<u>Step 2</u>:  $M = \{$ vertices with the least number of ancestors amongst those that have more ancestors than non-ancestors $\}$

Query on a vertex with the most balanced ratio amongst vertices of M or parents of vertices of M
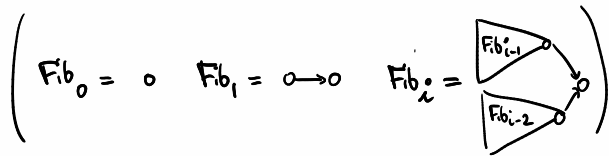
<u>Step 3:</u> Recurse

# THEORITICAL ANALYSIS OF THE GOLDEN BISECTION

> **Theorem**
>
> The *golden bisection* is a $\frac{1}{\log_2(\phi)}$ - approximation algorithm for binary DAGs, where $\phi$ = golden ratio
>
> $\frac{1}{\log_2(\phi)} \simeq 1,44$ is the optimal constant.

Problematic DAG for the *golden bisection* = $\text{Comb}(\text{Fib}_n)$ where $\text{Fib}_n$ is the $n$-th Fibonacci tree.



$$\left( \text{Fib}_0 = \circ \quad \text{Fib}_1 = \circ \rightarrow \circ \quad \text{Fib}_i = \right)$$

# PERSPECTIVES

→ Write an article

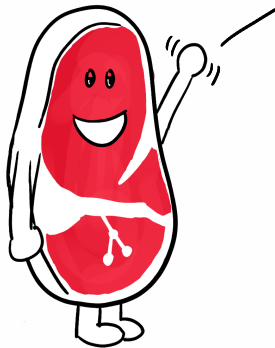→ Experimental results : git bisect Vs golden bisection

→ Average-case analysis
- Good model of random git graph (not Erdös-Rényi) ?
- How to sample them
- Theoritical analysis of git bisect

→ Efficiency of git bisect on trees ?
  Conjecture : git bisect = 2-approximation on trees