# How we measure and optimize for RAIL in V8's GC

mlippautz@, ulan@
BlinkOn 6, Munich

# Contents

**V8 GC and RAIL** → **Update on Optimizations** → **Reflection and Future**

Google

# Who We Are - V8 GC Team Munich

**Hannes**
Payer

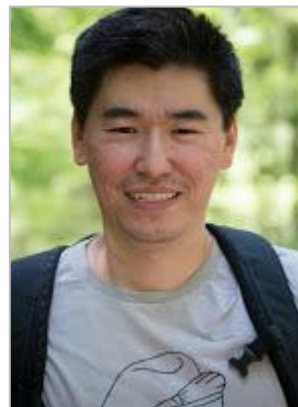**Jochen**
Eisinger

**Marcel**
Hlopko

**Michael**
Lippautz

**Ulan**
Degenbaev

# V8 GC and RAIL

# The Impossible Garbage Collection Triad

High Throughput

😔

Battery: Race to finish

Low Latency

Low Memory

# RAIL and GC Metrics

- Ideally: No GC... ever
- In practice: Prioritize

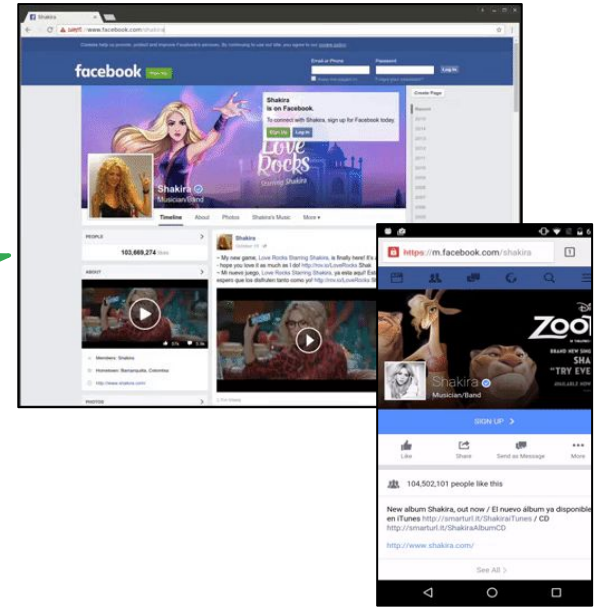| Response | Animation | Idle | Load |
|----------|-----------|------|------|

Max pause time          Overall GC time

*Fine print*

*Without consuming too much memory* 😊

Google

# Real World Benchmarking

- Catapult (Telemetry)
- Record/replay for real-world websites
- Mobile and desktop

*e.g. infinite scroll*

*e.g. initial page load*

?

R **A** I L

*idle website, e.g., gmail*



Google

# V8's Generational Garbage Collector

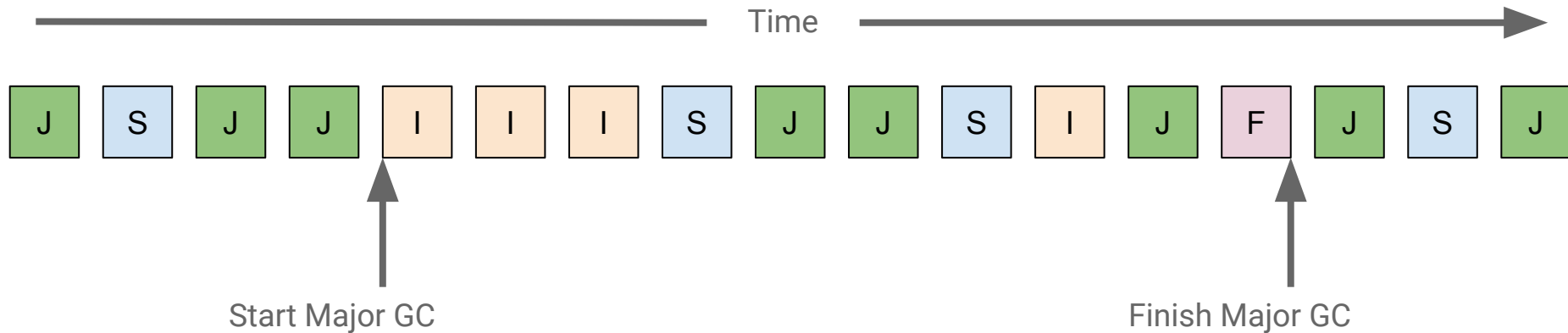Generational hypothesis: "Most objects die shortly after their allocation"

- New generation: Semi-space Scavenger (Cheney)
- Old generation: Mark-Compact (and Sweep)

New generation

Old generation

Page

Google

# GC Events

**J**     JavaScript code
**S**     Minor GC: Scavenger (~0-10 ms)
**I**      Major GC: Incremental Marking (~0.01-*CONFIGURABLE* ms)
**F**     Major GC: Final Mark-Compact Collection (~4-20 ms)

Time

J   S   J   J   I   I   I   S   J   J   S   I   J   F   J   S   J
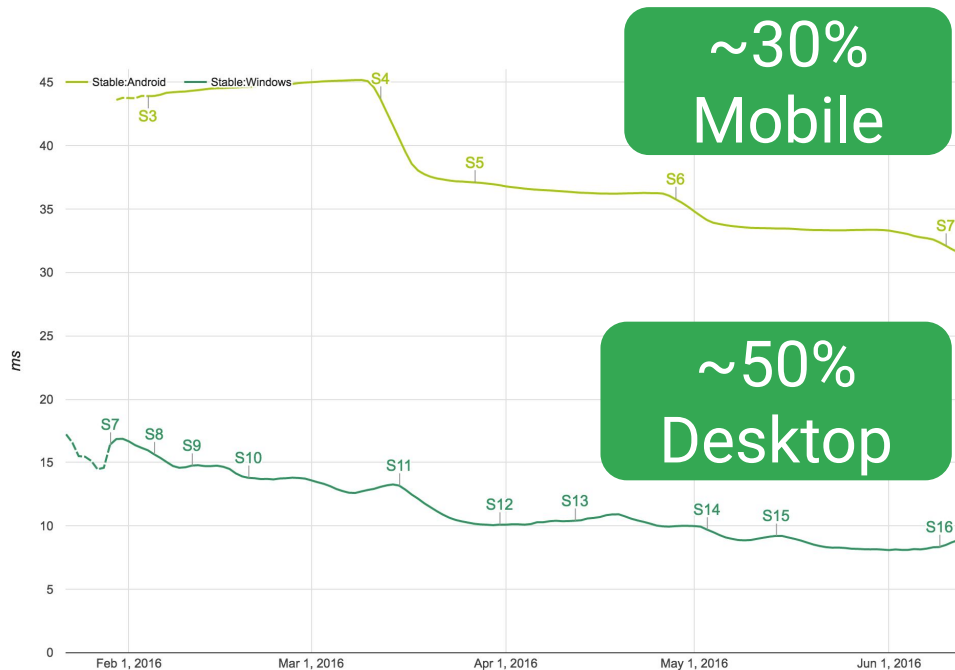
Start Major GC

Finish Major GC

Google

# V8 GC Update

Focusing on A

# Orinoco

Mostly parallel and concurrent GC
without strict generational boundaries
to
*reduce jank* and memory consumption
while providing high throughput

**... landing incrementally**



~30%
Mobile

~50%
Desktop

UMA: Major GC final pause **F** (50 %-ile)

Google

# Evacuation

- *Copy* objects within semi space of new generation
- *Move* objects from new to old generation
- *Copy* objects within the old generation
- Re-*write* remembered set pointers

**Writing memory is expensive**

*Orinoco*
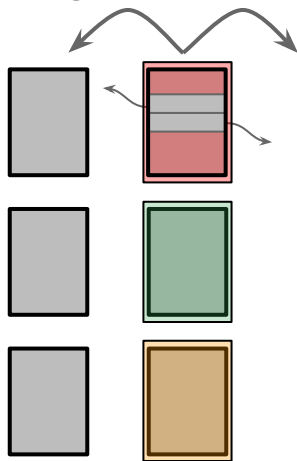
**Design GC to utilize available resources**
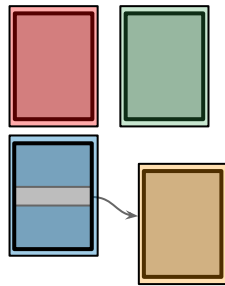
Google

# Parallel Evacuation during Major GC

- ## Lock-step
  - Parallelize moving of memory based on pages
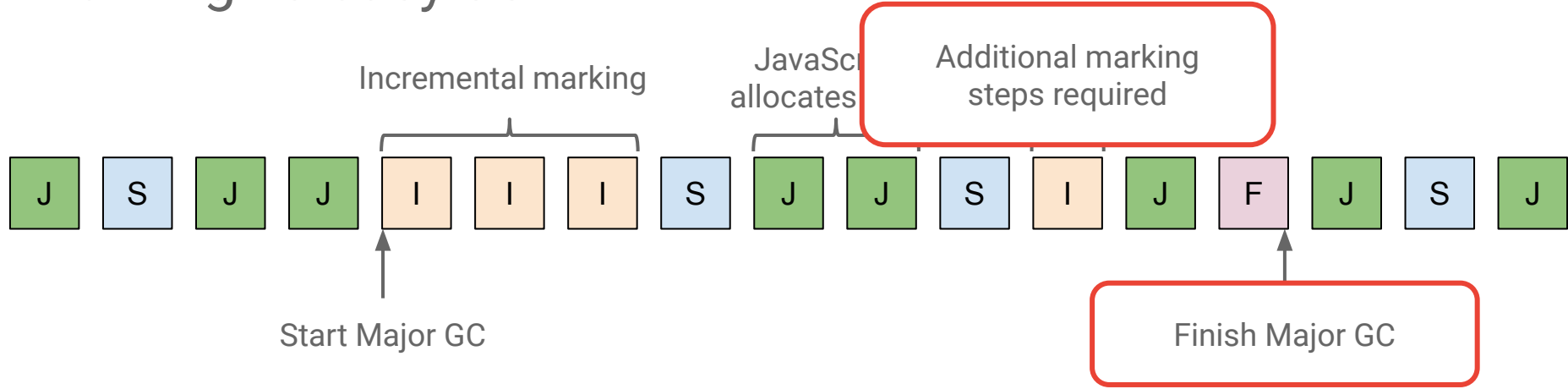  - Parallelize processing remembered set pointers
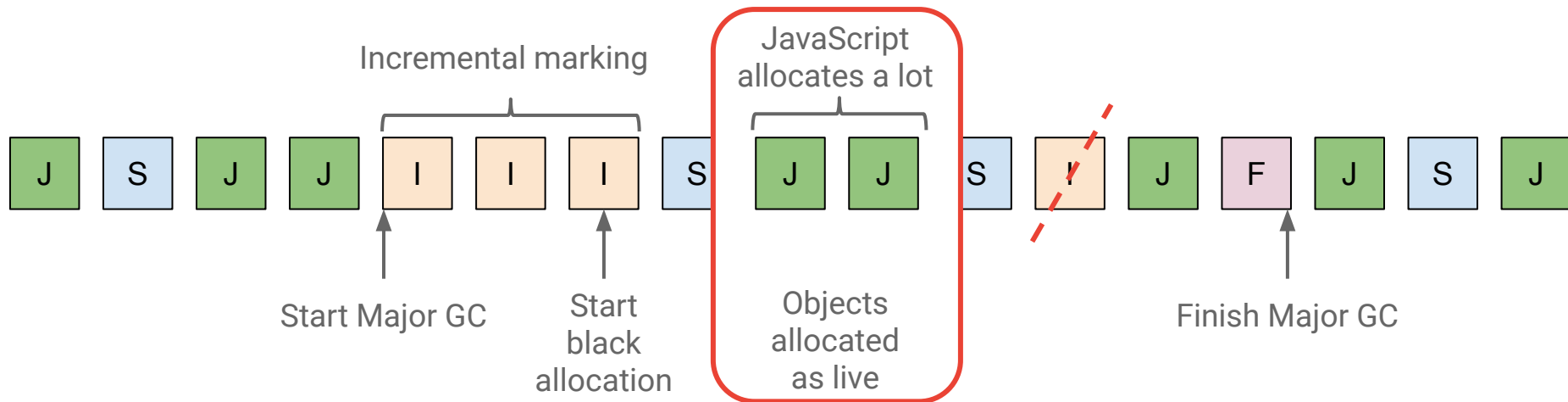
New generation

Old generation

Helper task

Main thread

Moving    Remembered set

# Marking vs busy JS

Incremental marking

JavaScript allocates

Additional marking steps required

| J | S | J | J | I | I | I | S | J | J | S | I | J | F | J | S | J |

Start Major GC

Finish Major GC

Live heap gets bigger
**More marking steps, longer finalization pause**

Google

🤘 Black Allocation

Incremental marking

JavaScript
allocates a lot

| J | S | J | J | I | I | I | S | J | J | S | I | J | F | J | S | J |

Start Major GC

Start
black
allocation

Objects
allocated
as live

Finish Major GC

Assumption: Objects allocated during marking will survive the following GC

Objects are allocated as live (already marked)
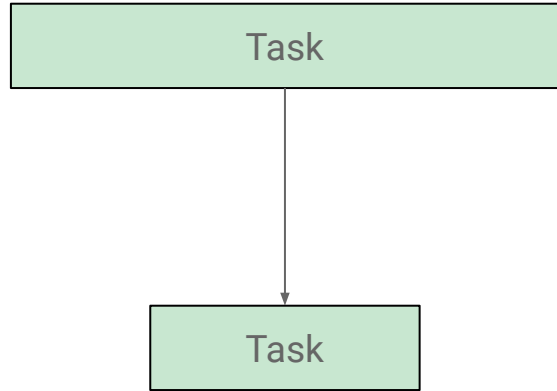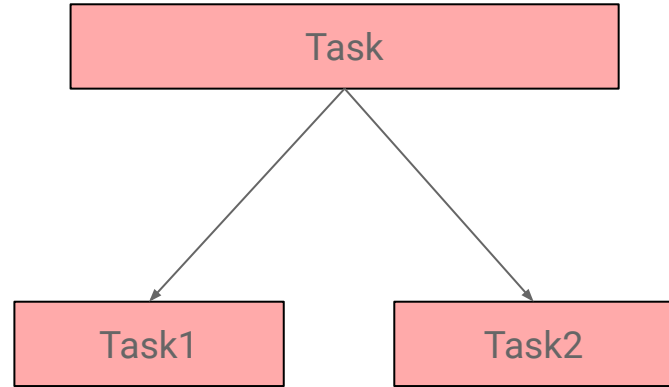
Google

# Some Results: Facebook during A



**Max mark-compact pause**

Time, ms (y-axis): 0, 7.5, 15, 22.5, 30

Chrome version (x-axis): 44, 45, 46, 47, 48, 49, 50, 51, 52

Google

# Some Results: Facebook during A

Average mark-compact pause



Google

# Reflection and Future

# Two ways to reduce latency

Optimize

Split up

| Task |
| --- |

| Task |
| --- |

| Task |
| --- |

| Task1 |
| --- |

| Task2 |
| --- |

Google

# Making GC more incremental

# Making GC more incremental

Before:

max-pause = 20ms

| Mark-Compact |
| :---: |

| A |
| :---: |

←————————— 20 ms ——————————→

After:

max-pause = 12ms

| A | Mark-Compact |
| :---: | :---: |

←— 8ms —→ ←———— 12 ms ————→

# Making GC more incremental

Before:

max-pause = 20ms

Mark-Compact

A

← 20 ms

**40% latency improvement?**

After:

max-pause = 12ms

A

Mark-Compact

← 8ms →

← 12 ms →

Google

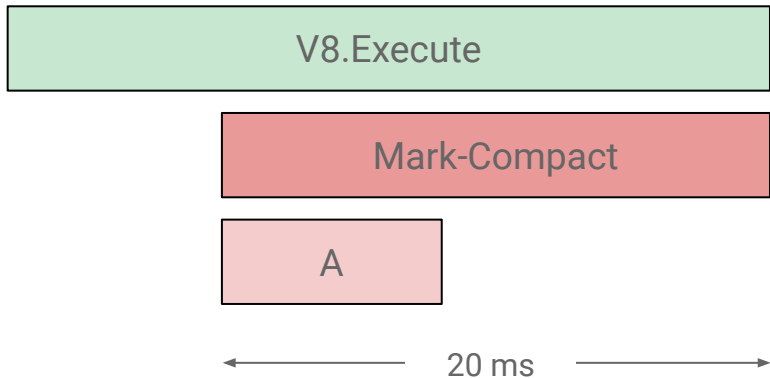# Making GC more incremental

Before:

max-pause = 20ms

| V8.Execute |
| --- |

| Mark-Compact |
| --- |

| A |
| --- |

← 20 ms →

After:

max-pause = 12ms

| V8.Execute |
| --- |

| A | | Mark-Compact |
| --- | --- | --- |

← 8ms → ← 12 ms →

# Making GC more incremental



Before:

max-pause = 20ms

V8.Execute

Mark-Compact

A

20 ms

After:

max-pause = 12ms

V8.Execute

A

Mark-Compact

12 ms

8ms

Google

# User facing metrics from PWM and system health

- **Response:** response latency
  - minimize GC work in critical time window
- **Animation:** animation latency
  - minimize GC work in critical time window
- **Idle:** responsiveness risk a.k.a responsiveness hazard
  - minimize GC work in tasks longer than 50ms
- **Load:** time until load finishes
  - minimize GC work during page load

Overfitting: delay GC in RAL, post many 49ms tasks in I

Simpler model: consider only RA, critical time window position is arbitrary

# Uniform incremental steps



Time Window

Time

Time Window

# Mutator Utilization



- Introduced by Cheng and Blelloch in 2001 for real-time GC.

- Called *mutator utilization* - application *mutates* reachability graph from GC POV

- Can be generalized to other tasks in Chrome
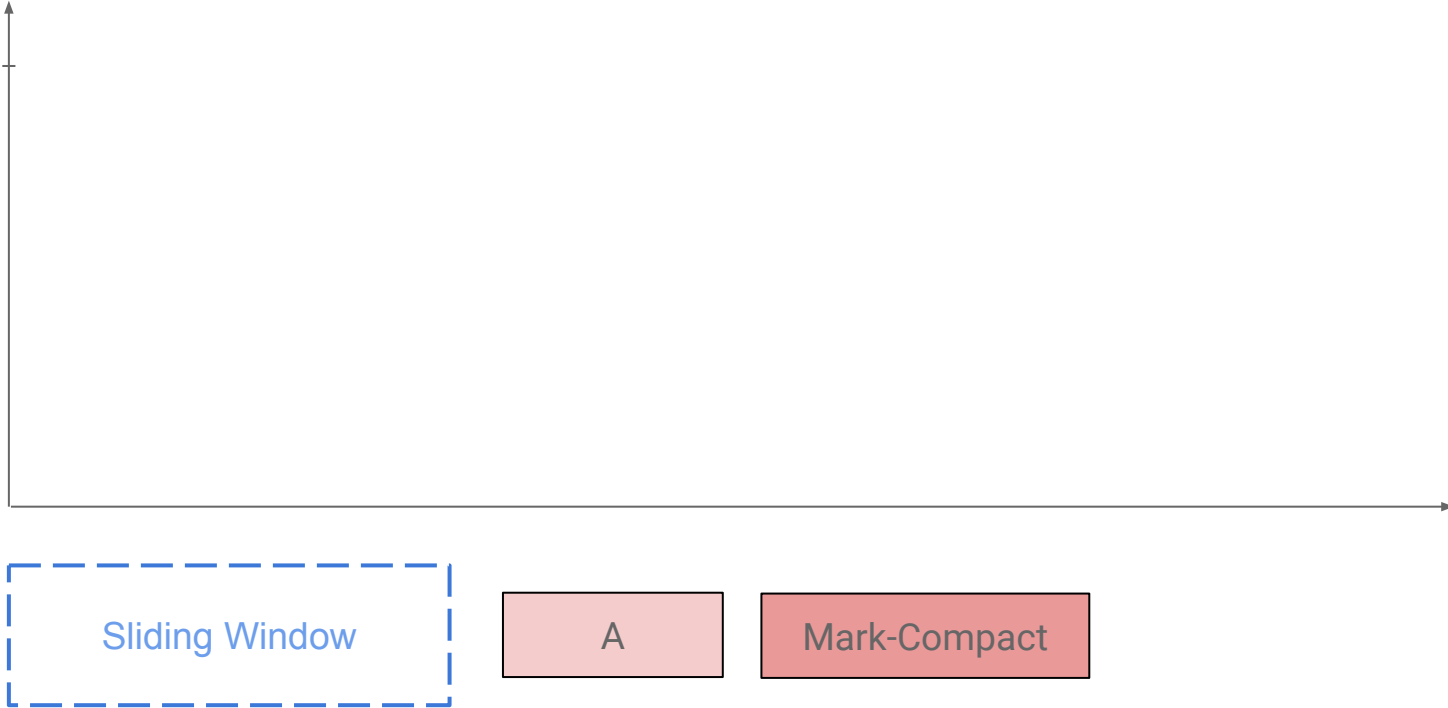
# Utilization function
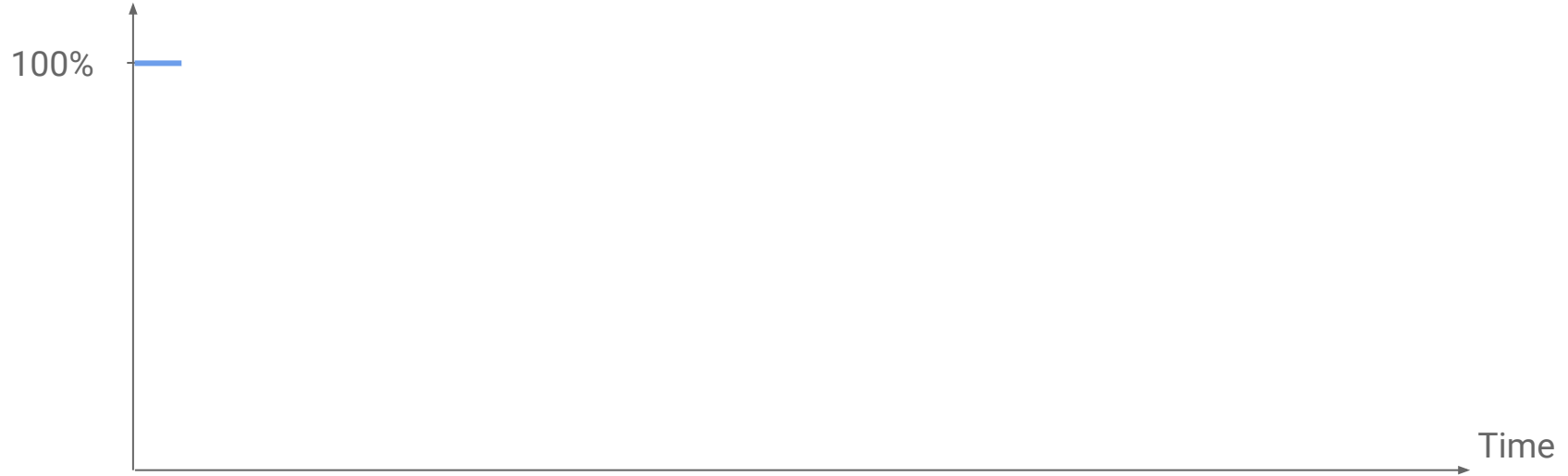
Free time in window

100% —

Time



Sliding Window

A

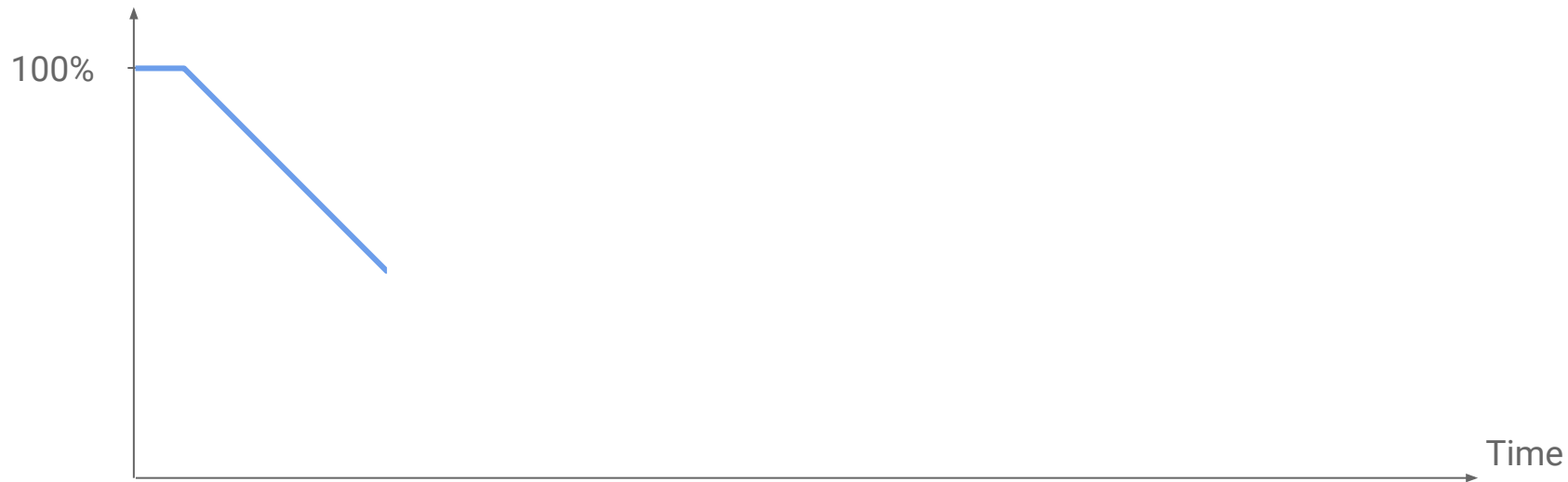Mark-Compact

# Utilization function

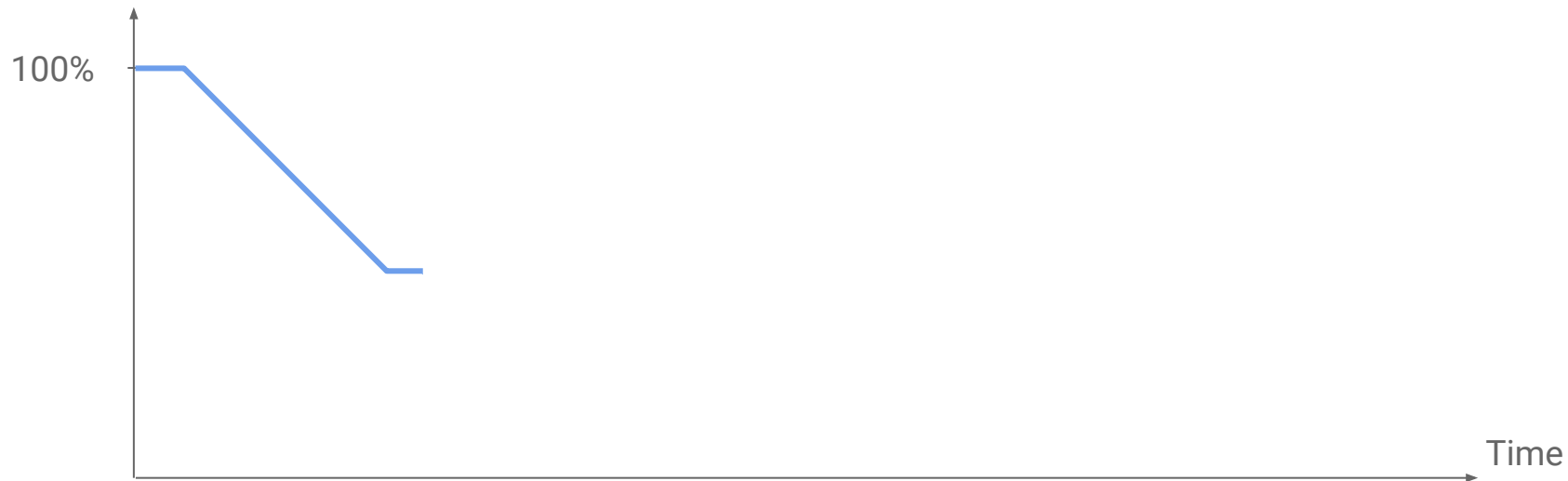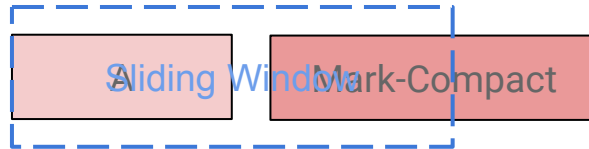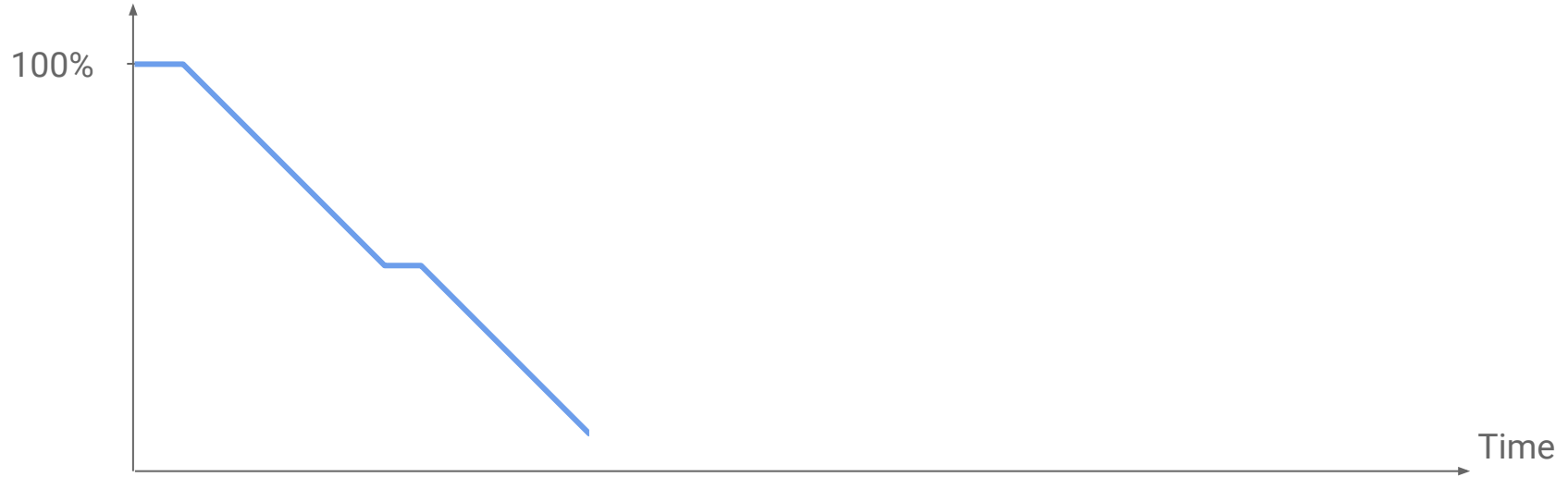Free time in window

100% —

Time

Sliding Window | A | Mark-Compact

# Utilization function

Free time in window

100%

Time

Sliding Window

Mark-Compact

Google

# Utilization function

Free time in window

100%

Time

Sliding Window

Mark-Compact

Google

# Utilization function



Free time in window

100%

Time

Sliding Window    Mark-Compact

Google

# Utilization function

Free time in window

100%

Time

A    Sliding-WindowMark-Compact

Google

# Utilization function

Free time in window



100%

Time

A          Mark-Compact

# Utilization function

Free time in window



100%

Time

A    Mark-Compact / Sliding Window

# Utilization function

# Utilization function

Free time in window
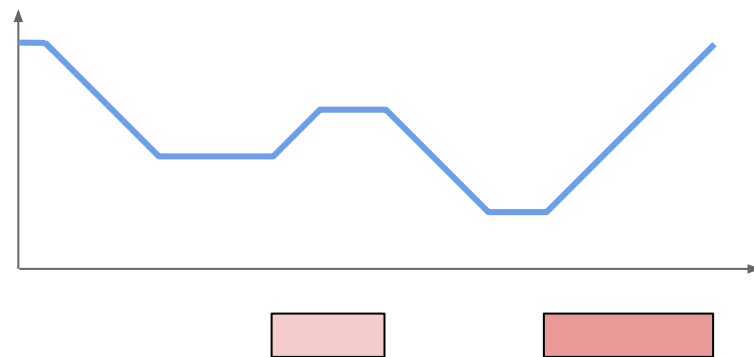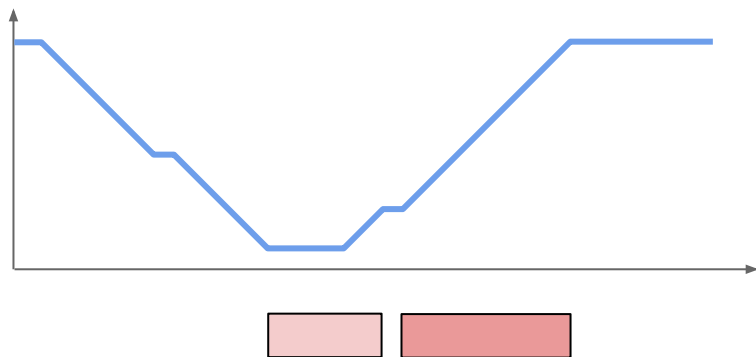
100%

Time

A

Mark-Compact

Sliding Window

Google

# Interpretations



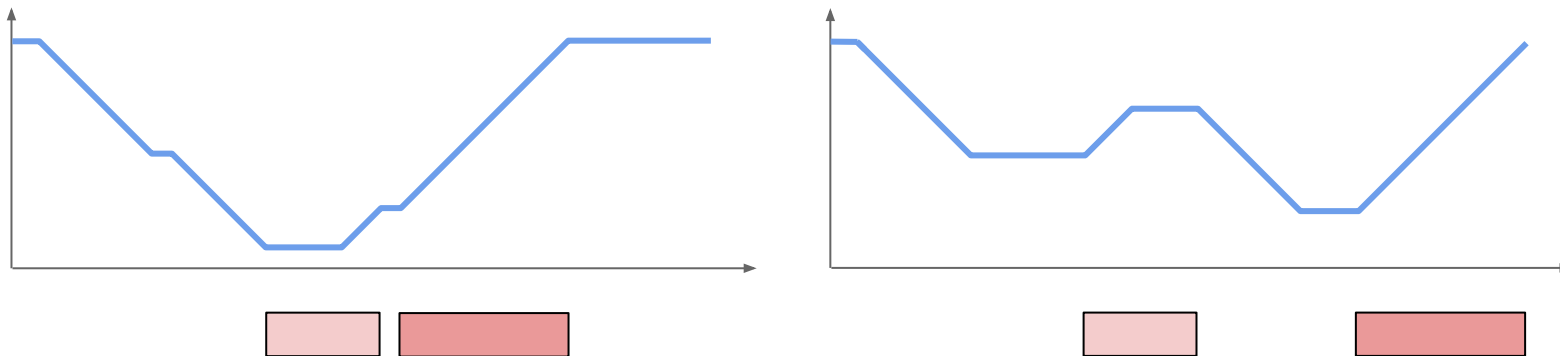- $f_w(t)$ = fraction of time GC leaves to application in the time window [t, t+w).

- $f_w(t)$ = probability that a high priority task arriving at any moment in [t, t+w) is not queued
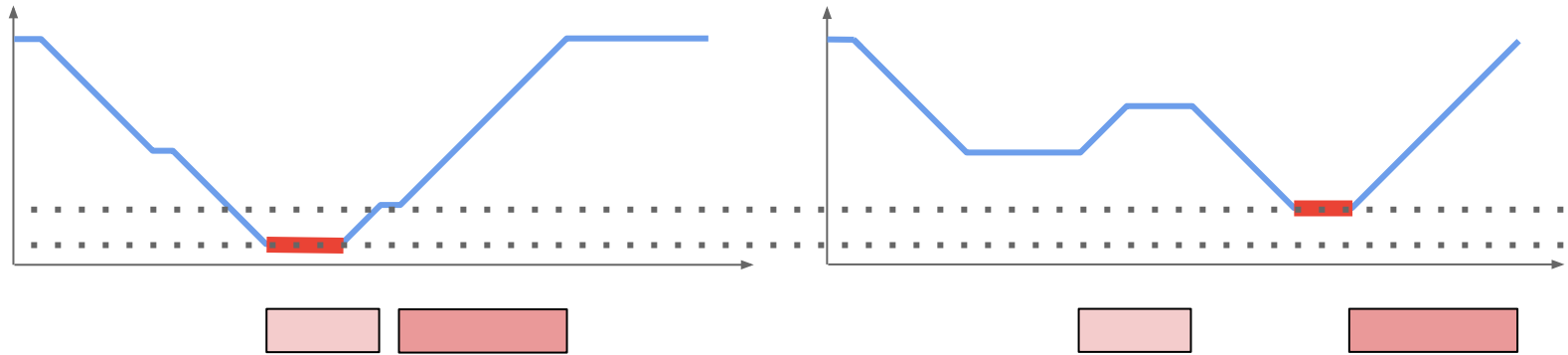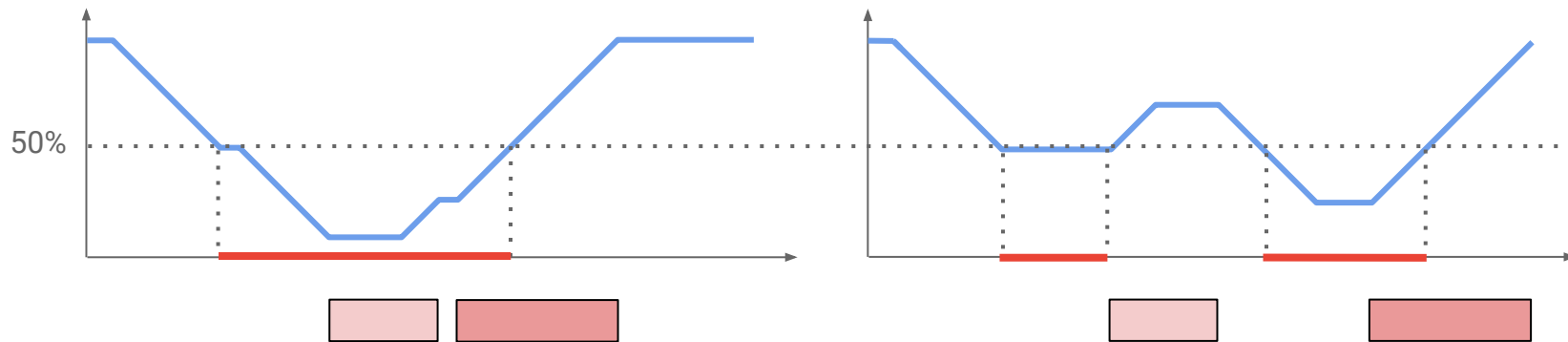
# Utilization function

# Utilization function



How to compare?
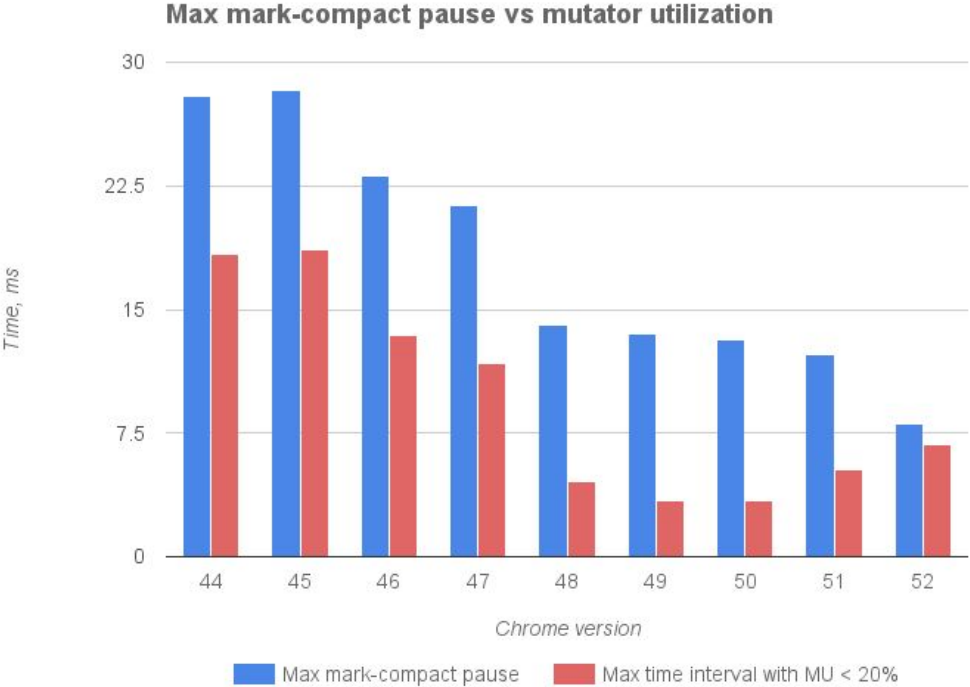
Google

# Minimum mutator utilization



Google

# Intervals with low mutator utilization



The largest interval where f(x) ≤ 50% for all points in the interval.

Google

# Evaluation



Max mark-compact pause vs mutator utilization

Google

# What we learned

- Incremental development of Orinoco - good idea!
  - Stability feedback via Canary channel
  - Performance feedback via UMA and telemetry
- Pause time distribution is not the whole performance story
  - Density of GC events on the timeline is also important
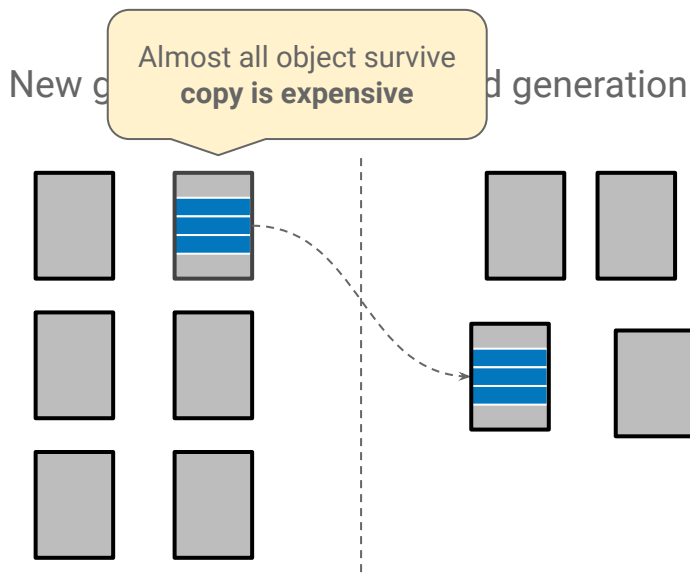  - Mutator utilization can be good indicator of latency impact

Google
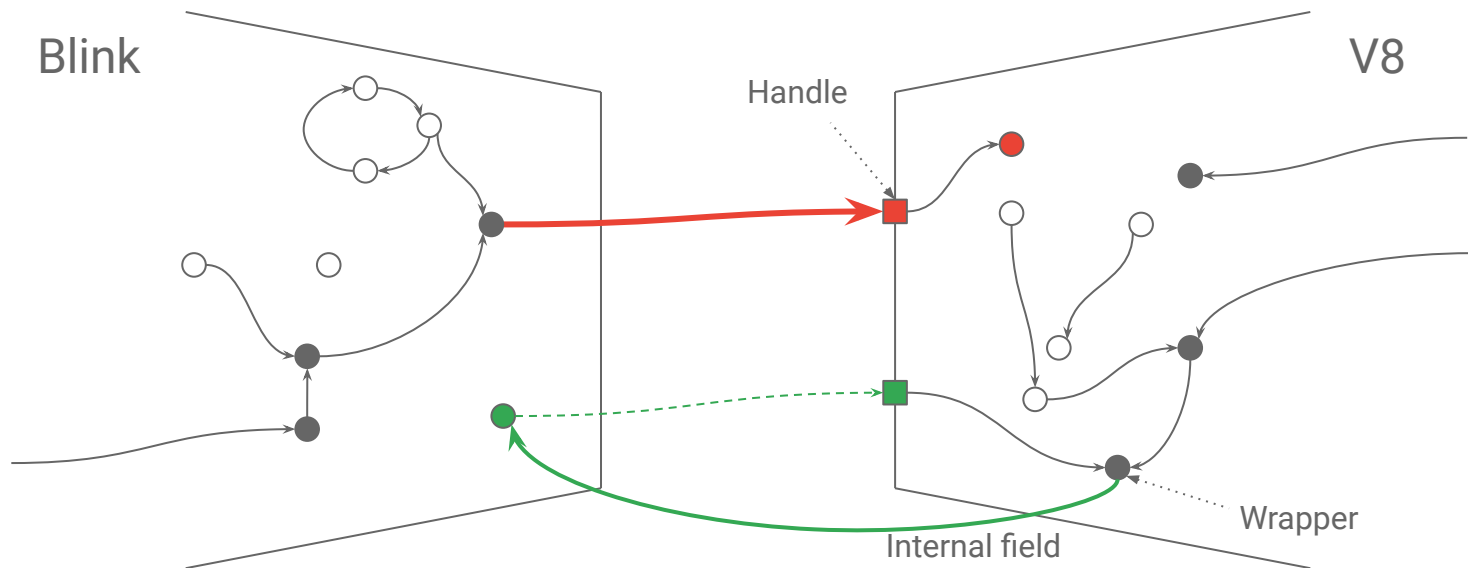
# Thanks!

# Changes on the way

## Experimental
(behind flags)

Google

# Unified Heap

- Allow pages to move from young to old generation without copying
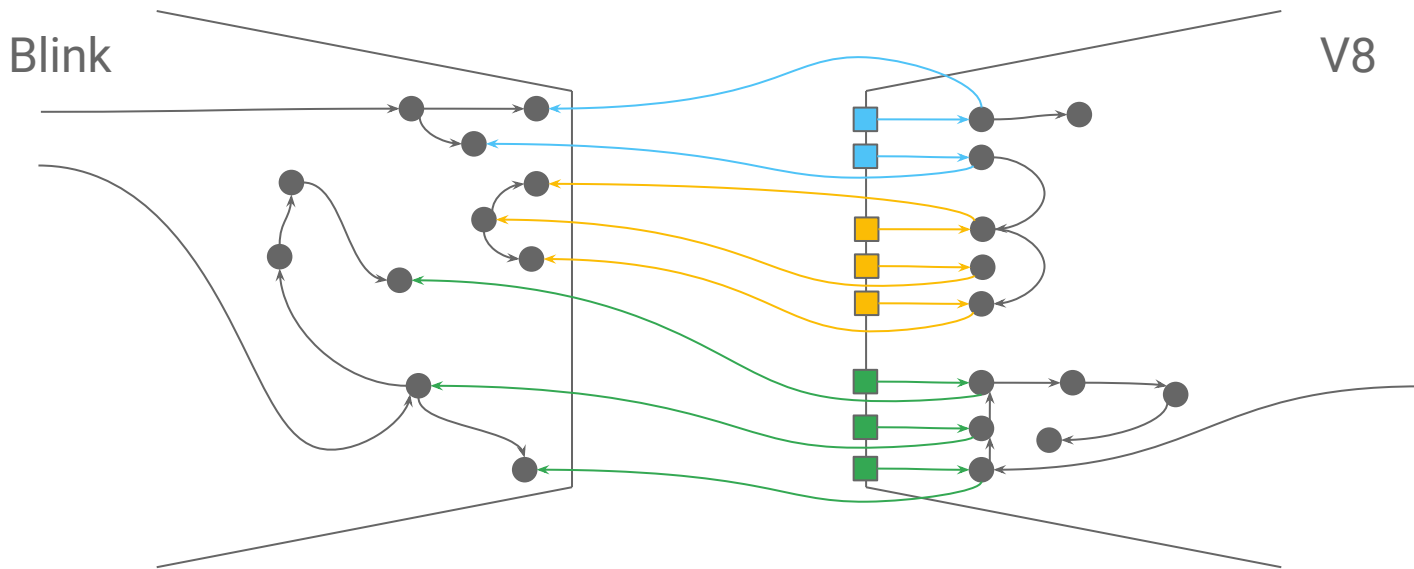
# Blink ⬌ V8 - Connecting the Worlds



Blink

V8

Handle

Wrapper

Internal field

Handles keep V8 objects alive
Wrapper keep Blink objects alive

Google

# Current Approach: Object Grouping

Blink

V8

- Handles know which group they belong too (rule-based)
- Groups keep objects alive

Google

# Becoming friends with Blink: Tracing of wrappers

Blink

V8

- Tracing through the Blink tree
- No need for handles

Google