

# Clojure Hash Maps:

plenty of room at the bottom

Peter Schuck  
@spinningtopsofdoom / @bendyworks

# Building an alien space ship

- Avoiding the gray goo scenario when making nano machines
- What cup of tea is best to power your Infinite Improbability Drive (earl gray hot)
- How to make the spaceship bigger on the inside than on the outside

Talk about real alien technology



# Optimizing Hash-Array Mapped Tries for Fast and Lean Immutable JVM Collections

by Michael J. Steindorfer and Jurgen J. Vinju

Compressed Hash-Array Mapped Prefix-tree

**CHAMP**

# ClojureScript Implementation

<https://github.com/bendyworks/lean-map>

# CHAMP gives you guaranteed Hash Map performance gains

- Iteration by 2x
- Equality checking by 10x to 100x



CHAMP trims your Hash Maps

going from relaxed fit punctuation

(

to their high school punctuation

{

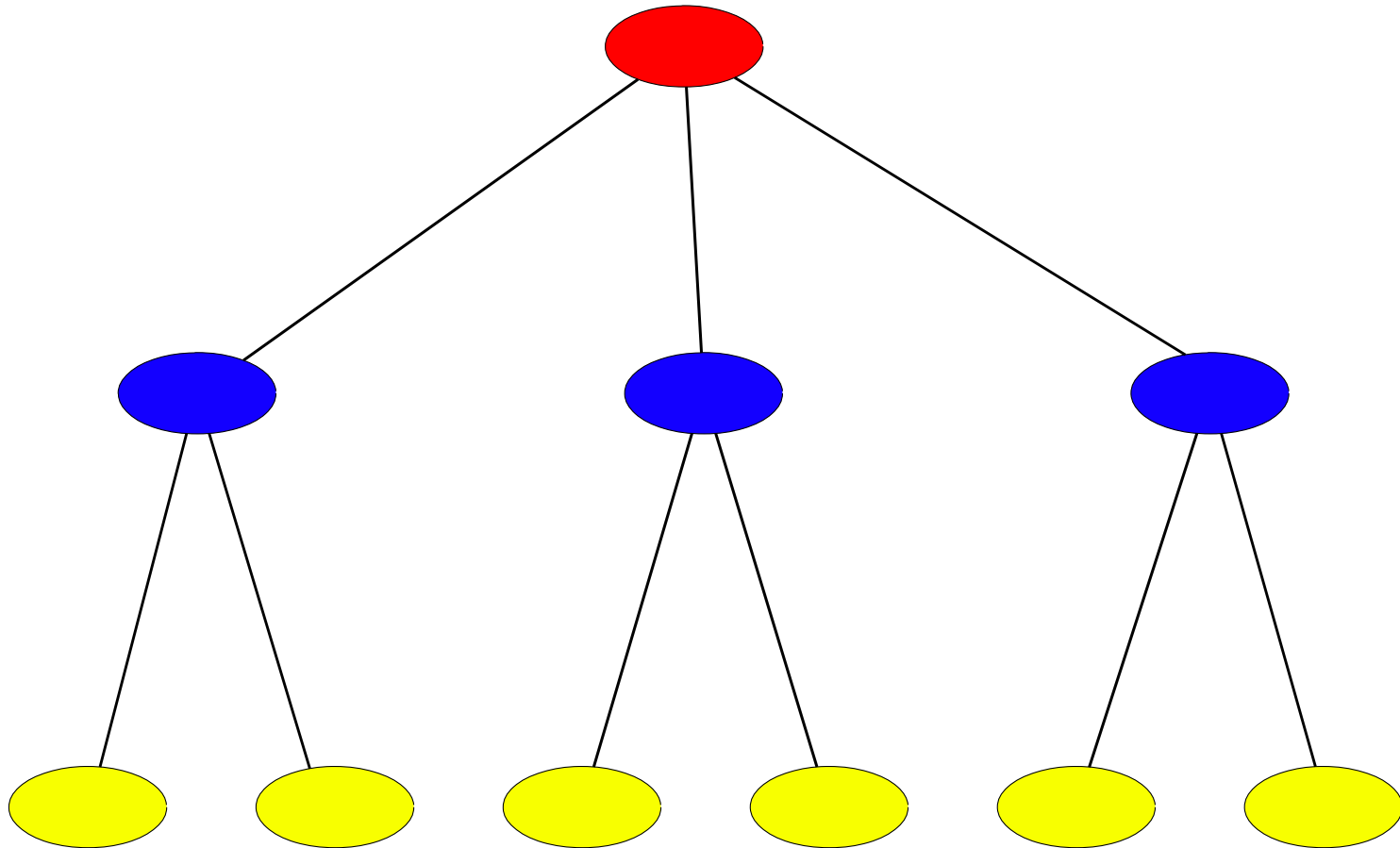


CHAMP makes Hash Maps more wieldy, making them both simpler and easier

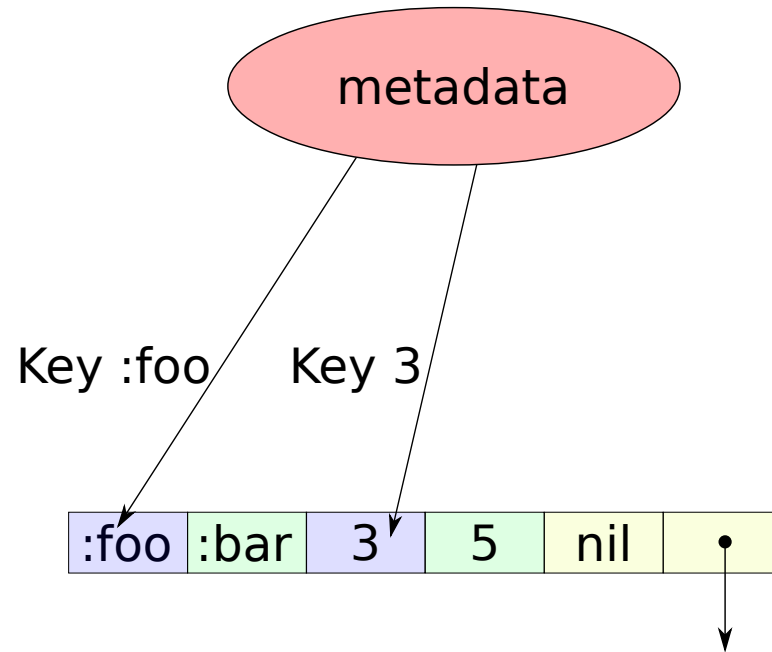
Code size is **two thirds** the size of the original implementation

# Overview of Clojure Hash Maps

# Clojure Hash Maps tree of nodes 32 way branching factor



# Node internals



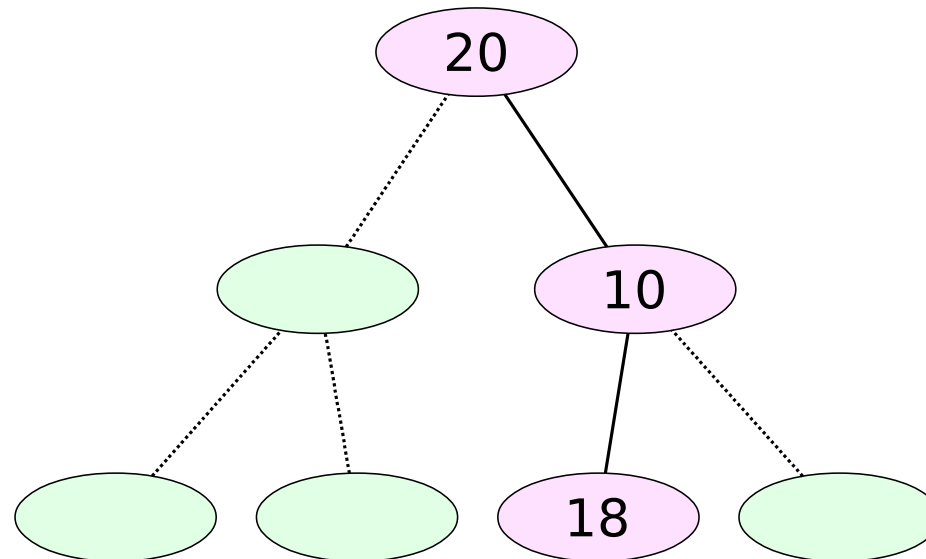


# How a key finds a node

Key: :foo

Hash: 1268894036

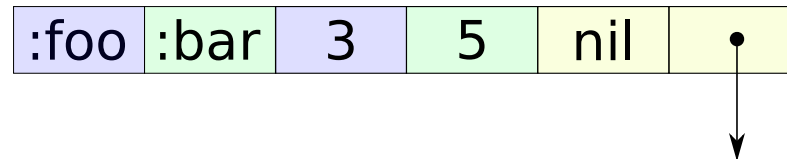
20	10	18	3	26	5	1
----	----	----	---	----	---	---



# First major improvement

Removes problems with sub node references

Sub node reference is a psuedo Key Value pair with `nil` as the "key"

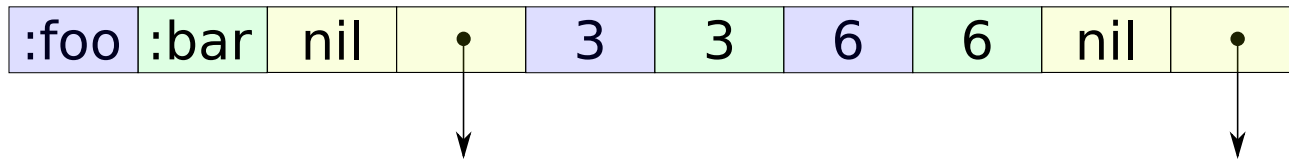


Doubles overhead for each sub  
node reference

# Adds incidental complexity

- Needs a flag for `nil` key and field for `nil` values
- Optimized node (Array Node) just containing sub node references
  - Happens when normal node's array has 32 elements
- Further complications with second problem

# Sub node references are scattered throughout a nodes array



Combined with `nil` marker value makes that for **every** operation you have to ask

"Is it a Key Value pair or sub node reference?"

Makes iteration a wiki walk



The Roman Empire was the post-  
**Roman Republic** period

The Roman Republic was the period  
of **ancient Roman civilization**  
beginning with the

Lots more link clicking...

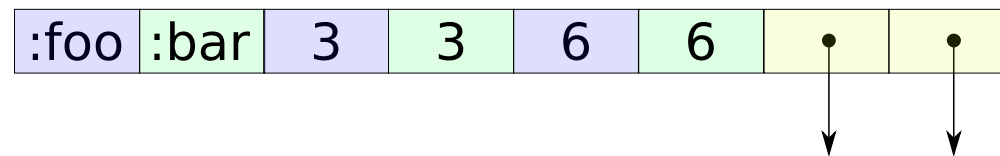
Awareness is the ability to perceive,  
to feel, or to be conscious of events,  
objects, thoughts, emotions, or  
sensory patterns

What was the next word after  
**Roman Republic?**

Wiki Walk Iteration has bad  
locality

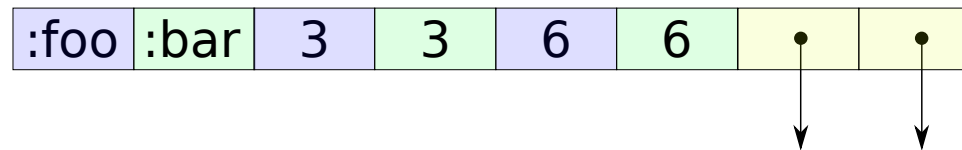
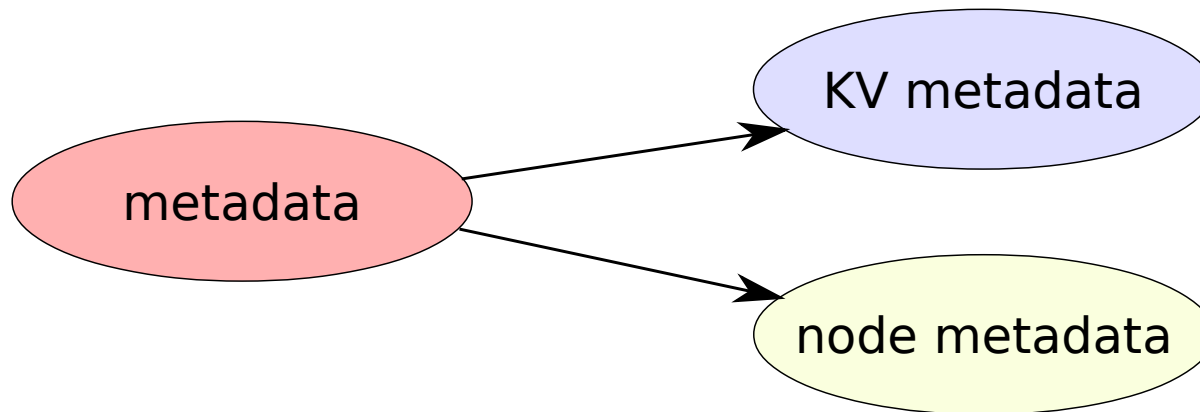
# CHAMP node improvements

# Key Value Pairs in front, Sub Node references in back





# Decomplect metadata



Lower memory overhead by  
removing `nil` marker values

# Removes all sub node incidental complexity

- nil key flag
- nil value field
- Array Node
- Check for Key Value or Sub node reference

2X performance by changing  
iteration from wiki walk to a linear  
scan

# Current Hash Map iteration algorithm

- If nil flag is true return [nil, <nil value>]
- For normal nodes
  - If key is not nil then return the Key Value pair
  - Otherwise go to sub node and repeat
- For Array node
  - If element is nil continue
  - Otherwise go to sub node and repeat

# CHAMP iteration algorithm

- Iterate through Key Value pairs
- Iterate through sub node(s) repeating step one

# Comparison

- Seven lines vs two lines
- Three conditionals vs none
- Polymorphism vs no polymorphism

# CHAMP Equality Check improvements



# Clojure Puzzler

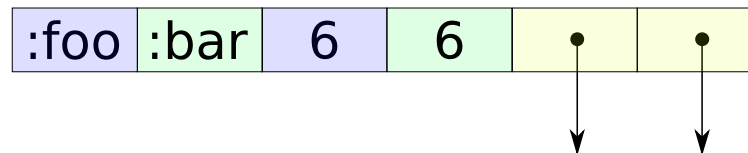
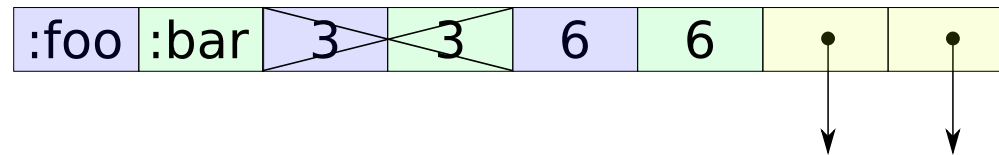
## Superficial Cleaning

```
(def base-map (hash-map)) ;; {}  
(def same-map  
  (as-> base-map m  
    (reduce #(assoc %1 %2 0) m (range 1000000))  
    (reduce #(dissoc %1 %2) m (range 1000000)))) ;; {}  
(= base-map same-map) ;; true  
(time (into {} base-map)) ;; 140 microseconds  
(time (into {} same-map)) ;; ??? microseconds
```

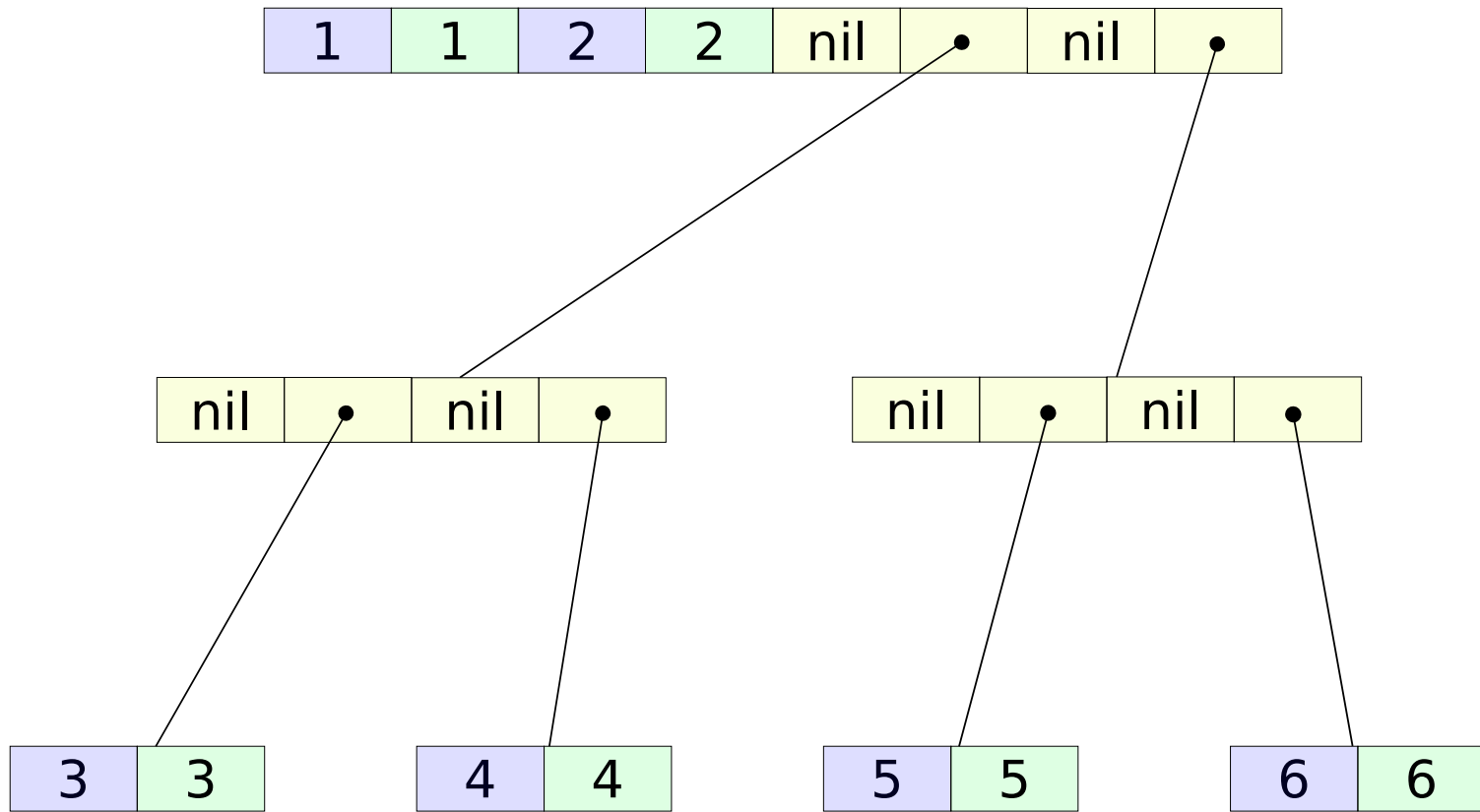
- A) 140 microseconds
- B) 280 microseconds
- C) 1500 microseconds
- D) 16500 microseconds
- E) 31000 microseconds

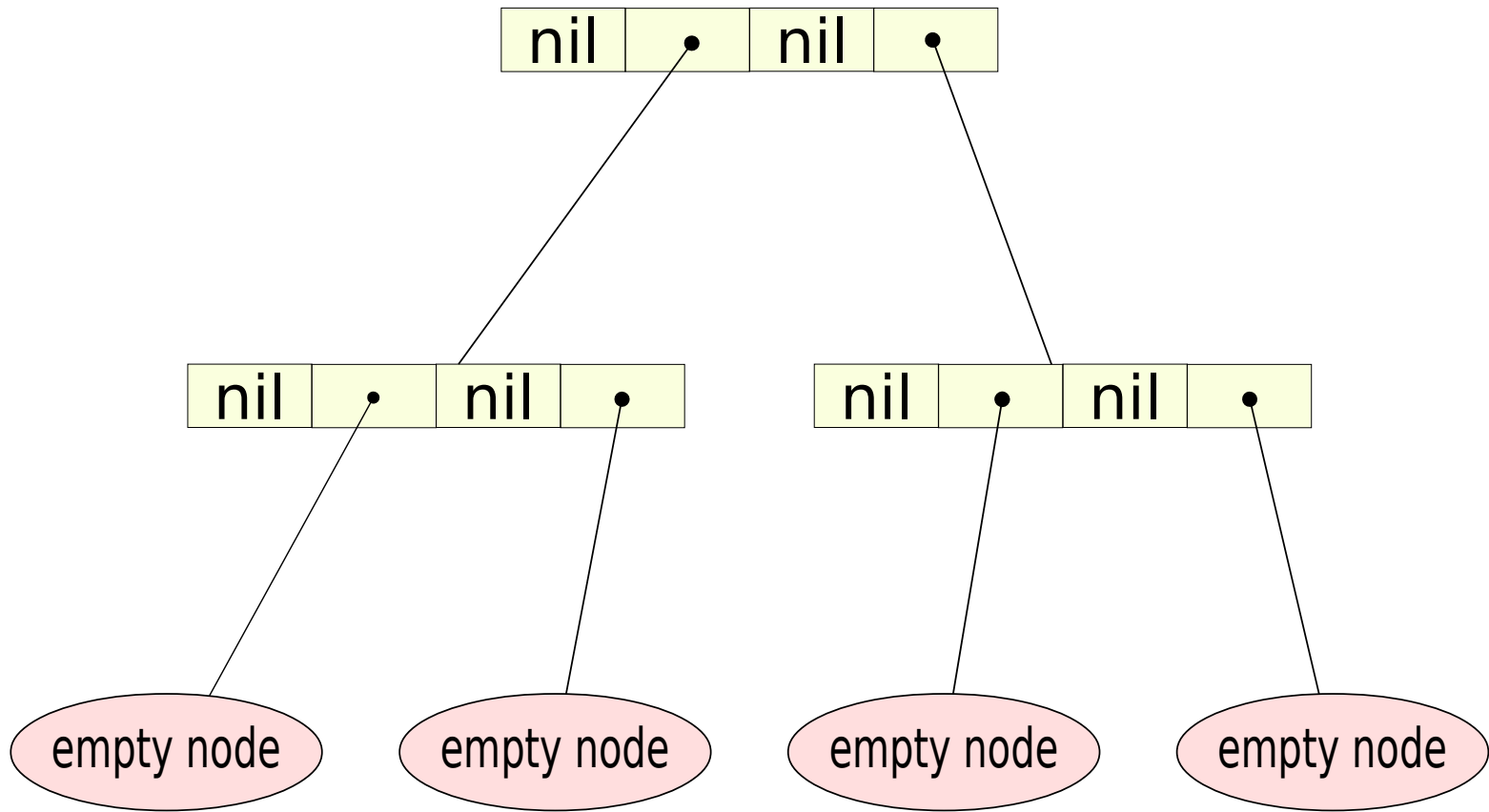
E) 31000 microseconds

# Current Delete Algorithm



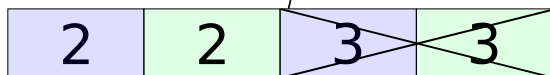
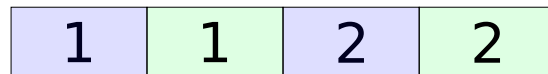
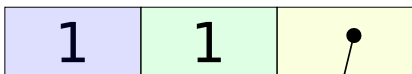
Superficial cleaning leads to

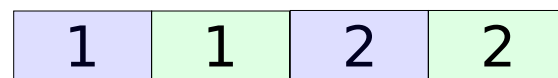
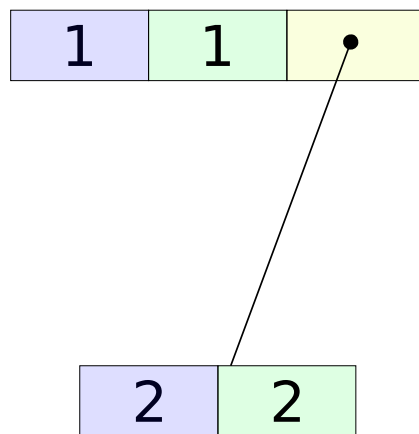
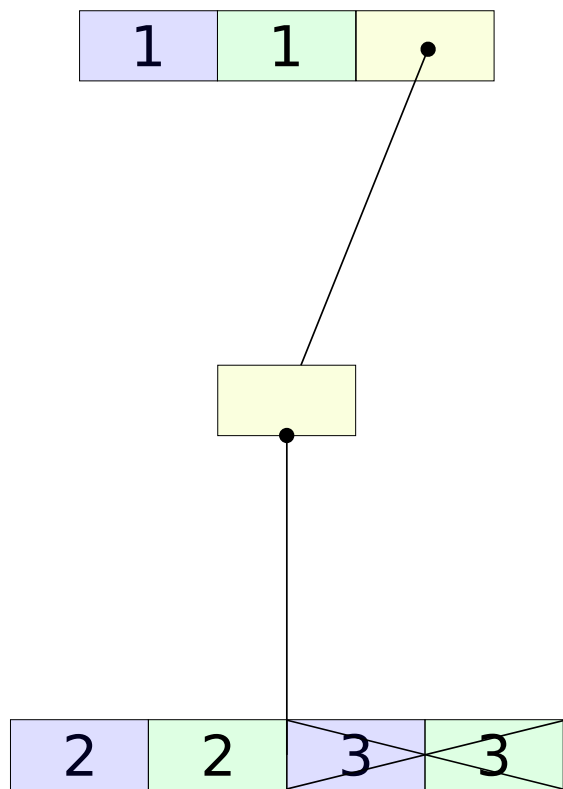






# CHAMP Delete Algorithm





Lowers memory overhead that  
occurs from `dissoc`

So what? This only really matters in pathological cases

So what? This only really matters in pathological cases

Equal CHAMP maps have the exact same layout in memory

So what? This only really matters in pathological cases

Equal CHAMP maps have the exact same layout in memory

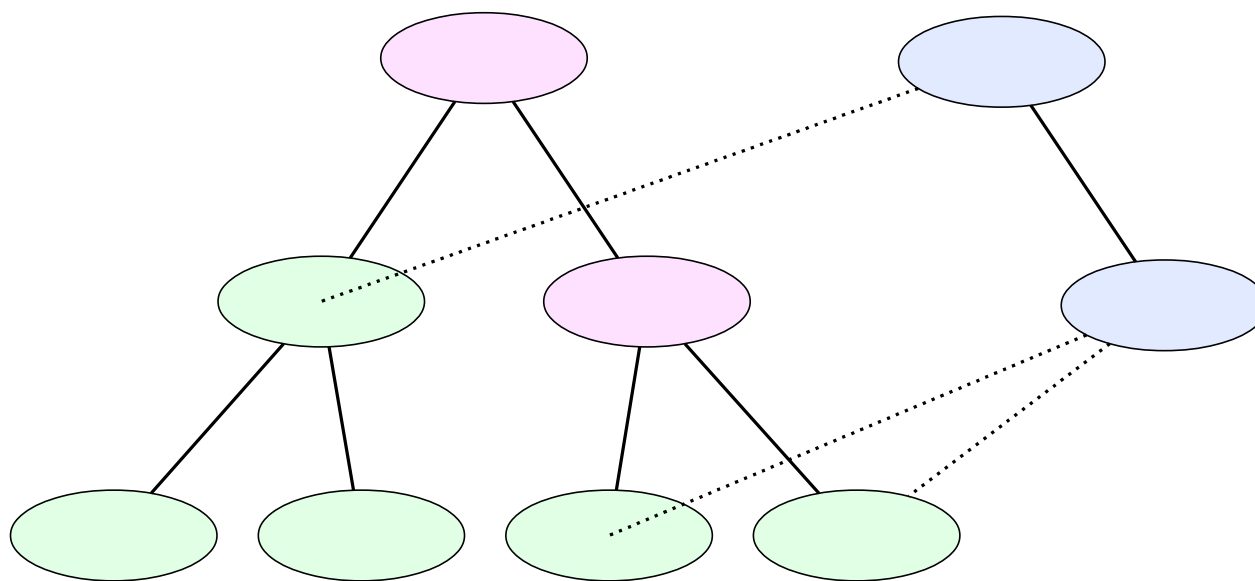
We don't have to compare all Key Values we can compare nodes

Equality check is now  $O(\log n)$  vs  $O(n)$  leading to 100x performance improvement

This is when maps share structure



# Structural Sharing



# We still get 10x performance boost for maps don't share any structure

- Current comparison has overhead due to Clojure abstractions (sequences and lookup)
- CHAMP comparison is only comparing two arrays

CHAMP improvements paves the  
way for future improvements

# Two Future possibilities

- Merge and Diff operations could have greatly increased performance
- Similar to RRB Vectors for Vectors

CHAMP is not as cool as working  
with nanobots



# CHAMP shows Hash Maps have plenty of room at the bottom

- 2x performance for iteration
- 10 - 100x performance for equality checking
- Lower memory overhead

For me biggest win is making Hash  
Maps much easier to understand  
and implement



# Clojure Hash Maps is one of Clojure's best exports

- Base Hash Map
  - Scala
  - Elixir
- Ports
  - Ruby (Hamster)
  - JavaScript (immutable.js)

# Thanks

- Bendyworks for supporting my work on this
- Michael J. Steindorfer and Jurgen J. Vinju for the CHAMP Paper
- Zach Tellman for writing Collection Check
- Martin Klepsch for porting Collection Check to ClojureScript
- Nicolás Berger for helping me setup test harness
- David Nolen for performance and profiling suggestions

Special Thanks

Cliff Rodgers

@2kliph

*Fin*

Questions?