

Rabin Fingerprinting The Hash Cannon



Arnaud Bellizzi
@Oodrive
Architect - Backup & Archive

Data !

A

B

C

D

E

F

G

H

Chunks !



Chunks !



- Differential transfer
- Deduplication

Byte Shift

A B C D E F G H

A B C Z D E F G H

Byte Shift

A B C D E F G H

A B C Z D E F G H

Byte Shift



CHANGED

CHANGED



What would we like ?

A B C D E F G H

A B C Z D E F G H

What would we like ?

A B C D E F G H

A B C Z D E F G H

What would we like ?

A B C D E F G H

CHANGED

A B C Z D E F G H

What would we like ?

- Variable chunk size
- Content based chunking

Enter Rabin !



Michael O. Rabin
1931

- Turing Award Winner 1976
- Great smile
- Loves Math, Cryptography, Hashes

Fingerprinting

- Hashing function (arbitrary lengths, statistically uniform)

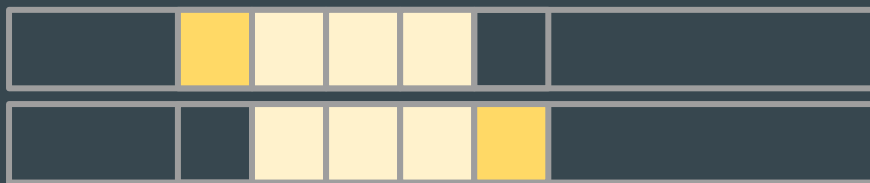


Fingerprinting

- Hashing function (arbitrary lengths, statistically uniform)



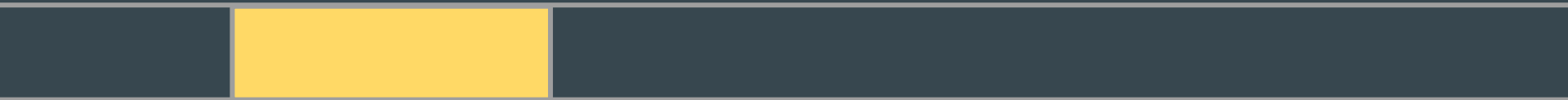
- Rolling hash



The Hash Cannon !



The Hash Cannon !

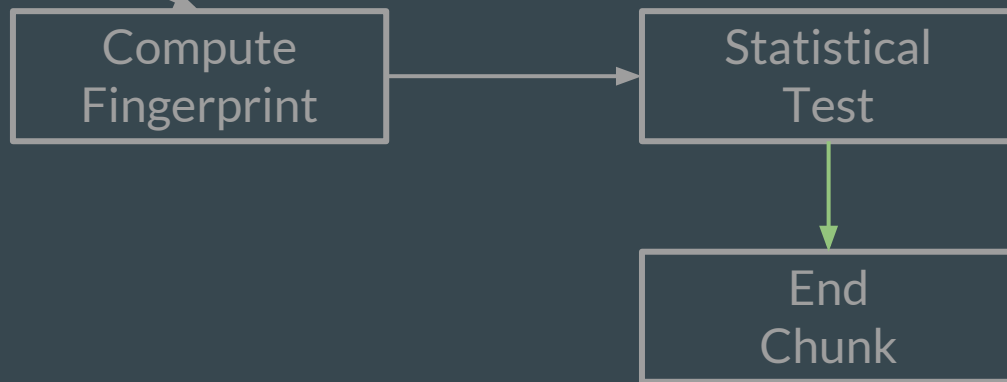


Compute
Fingerprint

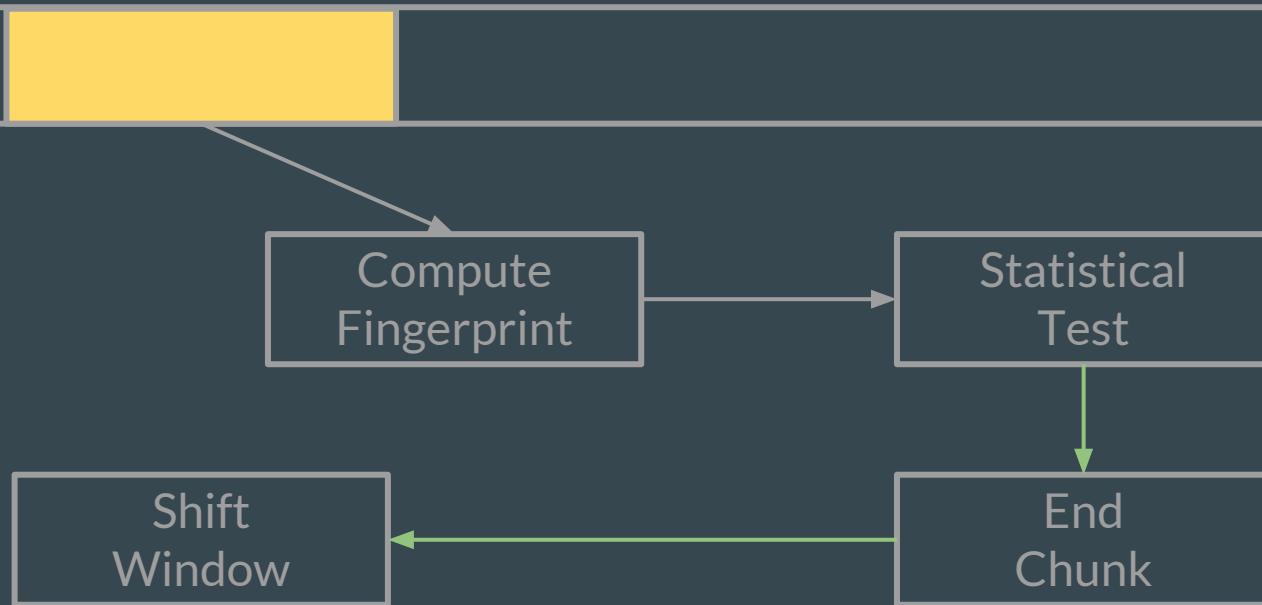
The Hash Cannon !



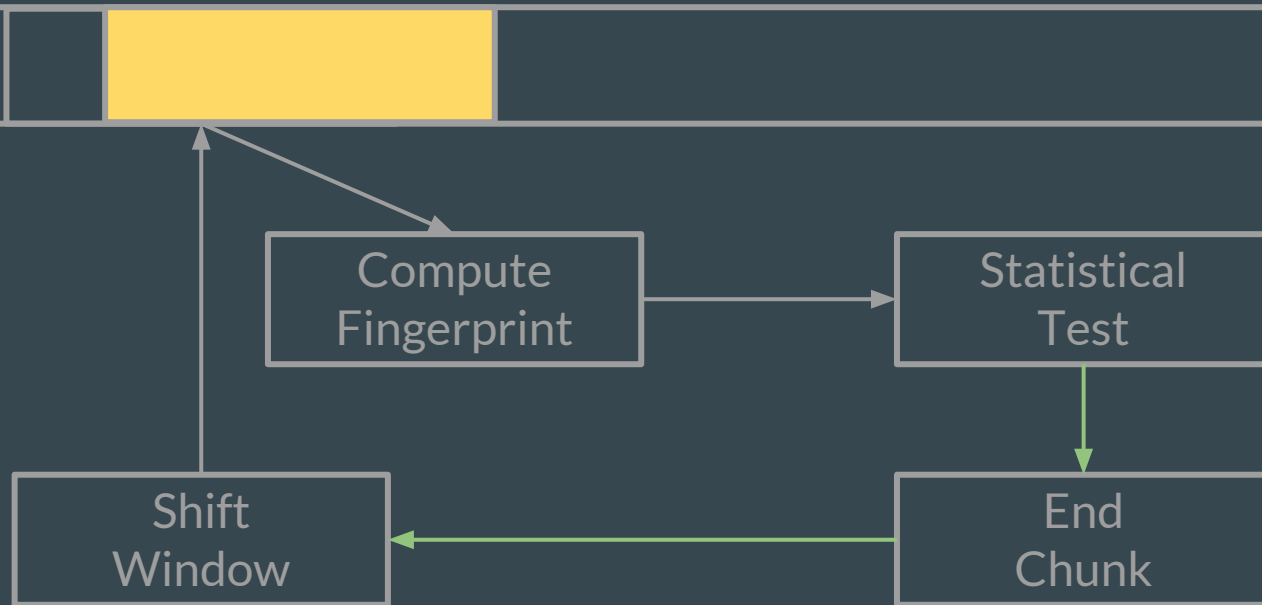
The Hash Cannon !



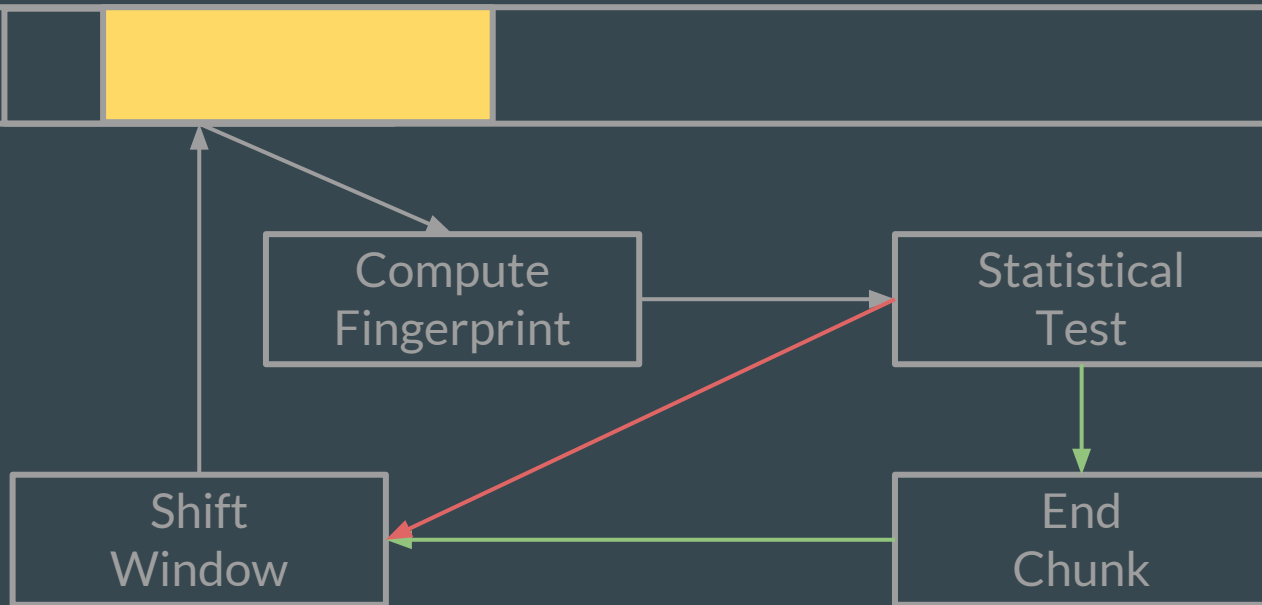
The Hash Cannon !



The Hash Cannon !



The Hash Cannon !



Byte Shift (again)

A B C D E F G H



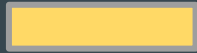
A B C Z D E F G H

Byte Shift (again)



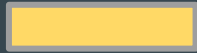
Byte Shift (again)

A B C D E F G H

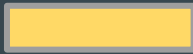


A B C Z D E F G H

Byte Shift (again)



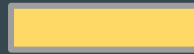
Byte Shift (worst case)



Byte Shift (worst case)



Byte Shift (worst case)



Byte Shift (worst case)



Byte Shift (worst case)



Byte Shift (worst case)



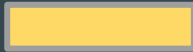
Byte Shift (worst case)



Byte Shift (worst case)



+1 CHANGED
+1 CHUNK



Profit !

~20% less chunks for same data

3-15% less data on the wire

Chunks are now **resistant** to Byte Shift

Thank you