

Add ALL
the things

avi@stripe.com





Steven H. Noble
@snoble



Following

If you are building a system for calculating aggregates and you don't know the relevance of things like abelian groups... STOP!!



Reply



Retweet



Favorite



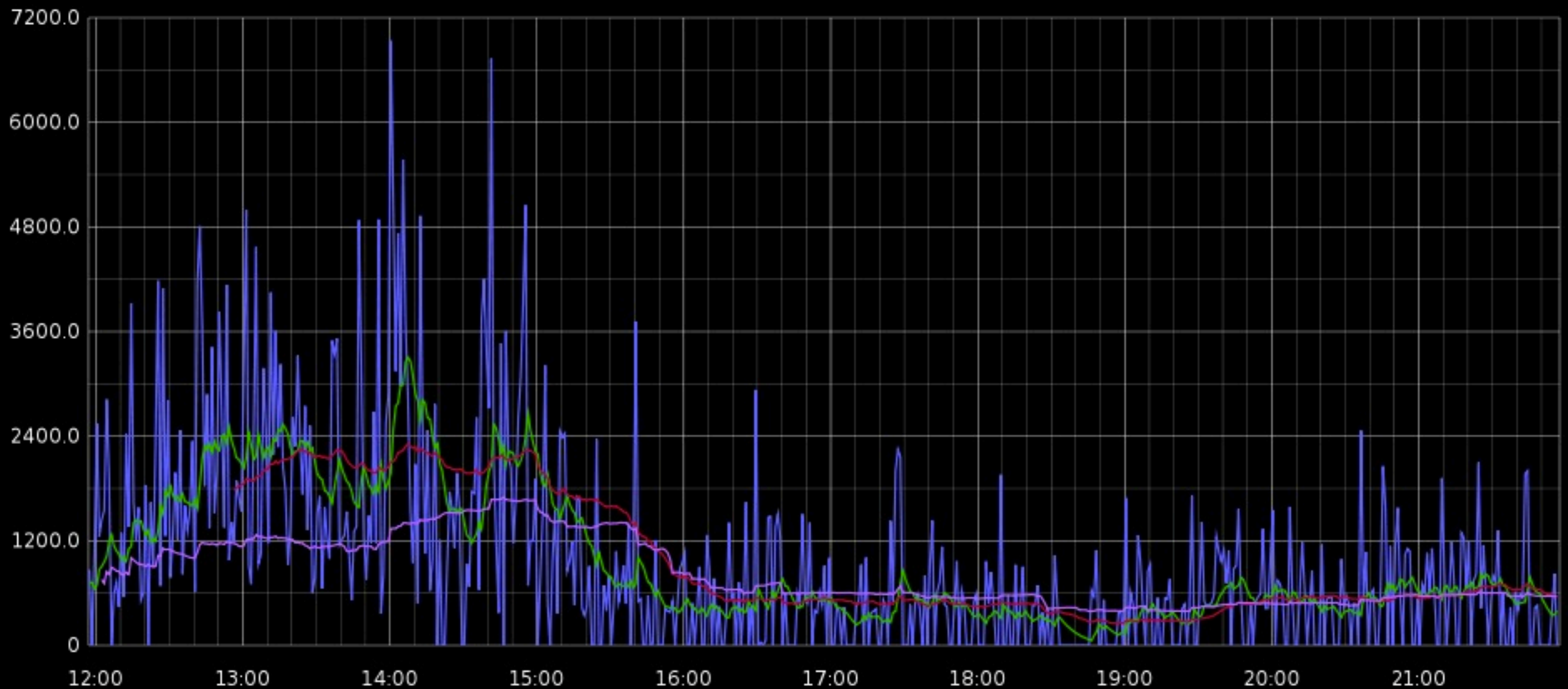
More

lolwut

tl;dr

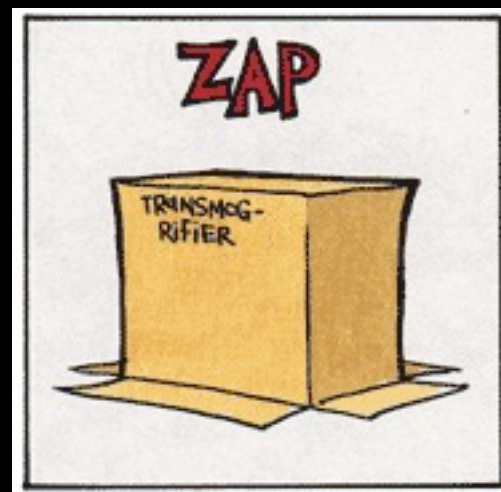
- Adding is awesome
- A **lot** of things that aren't adding are still "adding" (which is awesome)

Motivating example: StatsD-like



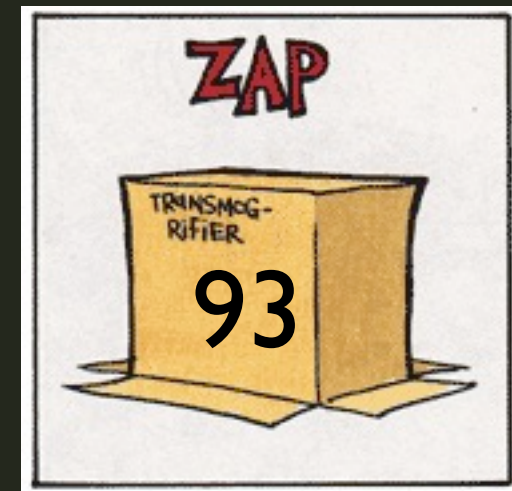
```
1 module Stripe
2   def process_payment(pmt)
3     stats.increment("money", pmt.amount)
4     transfer(pmt.from, pmt.to, pmt.amount)
5   end
6 end
```

Addifier

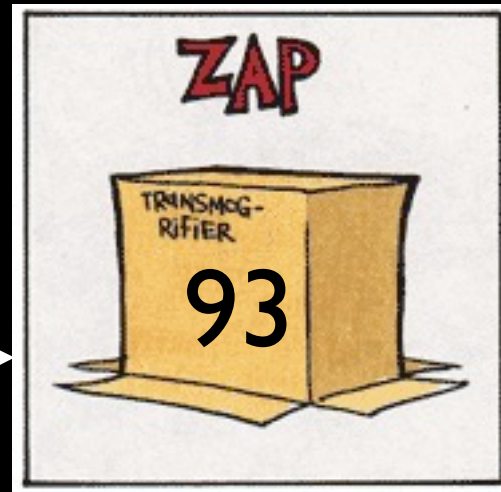


4 + 8 + 15 + 16 + 23 + 42

```
8 class Addifier
9   def initialize
10     @result = 0
11   end
12
13   def update(x)
14     @result = @result + x
15   end
16 end
```



4 8 16 23 42



$$((((4+8)+15)+16)+23)+42$$

12:00

12:01

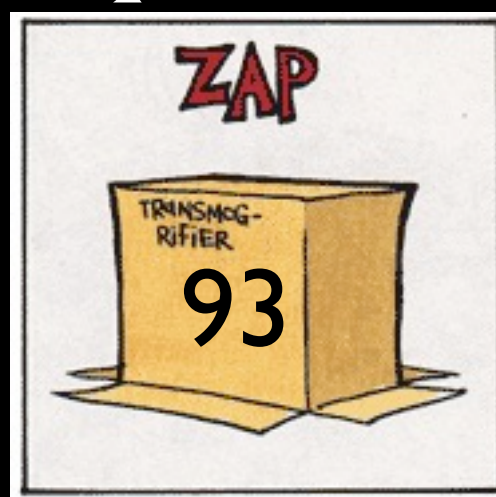
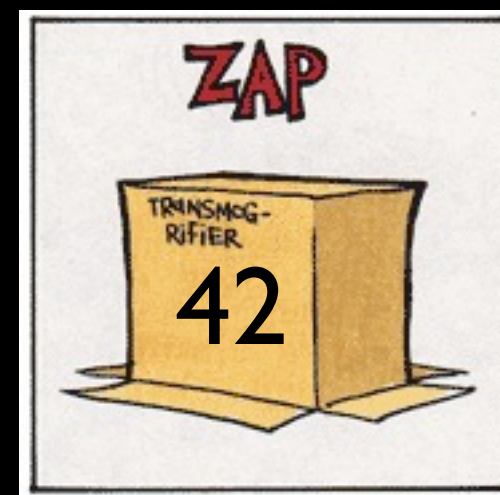
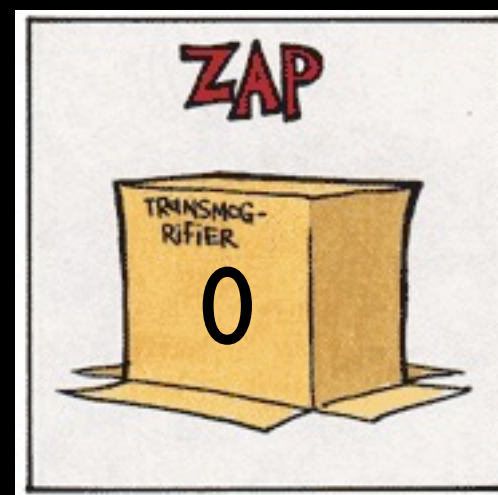
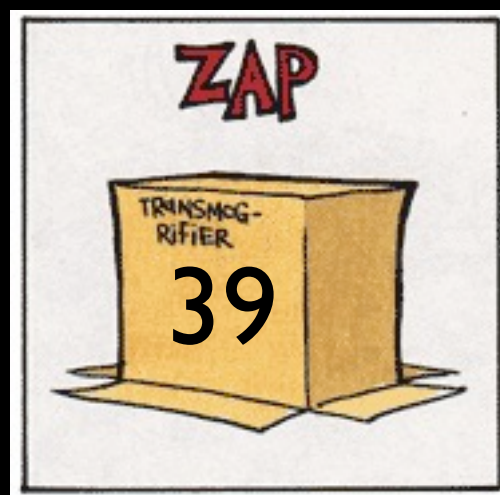
12:02

12:03

4 8

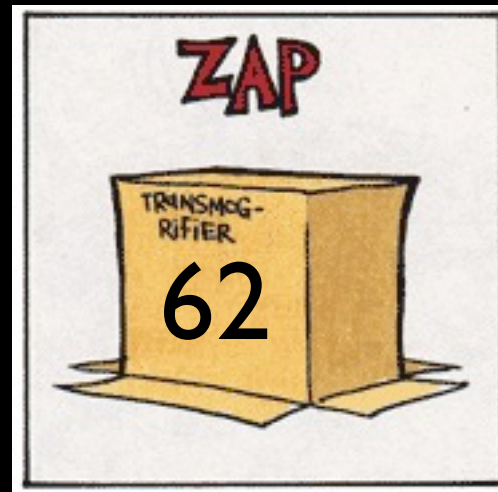
16 23

42

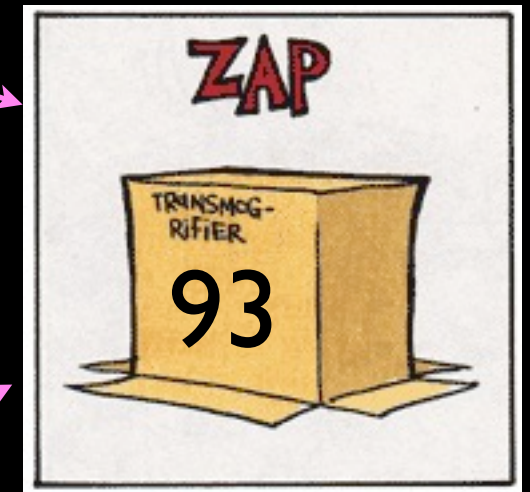
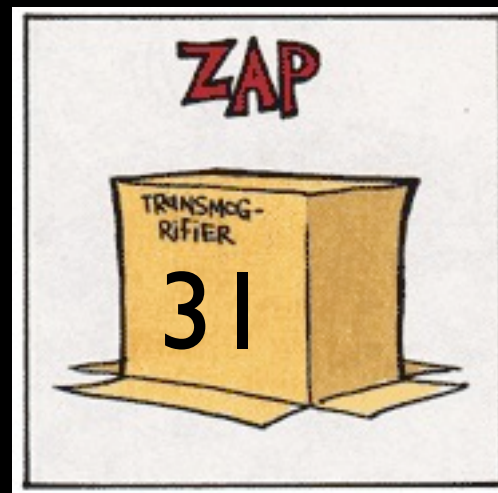


$$(4+8)+(16+23)+0+42$$

4 16 42
→



8 23
→



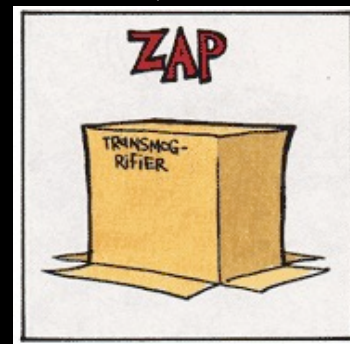
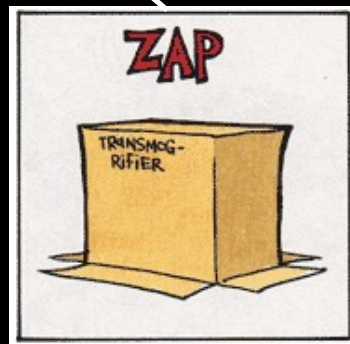
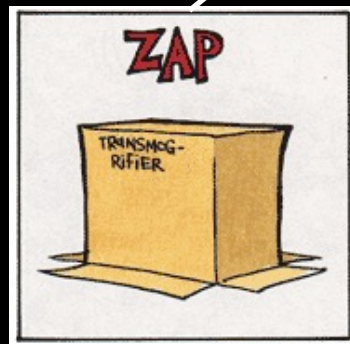
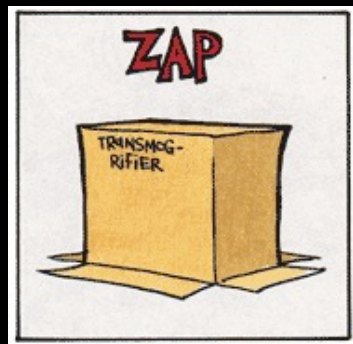
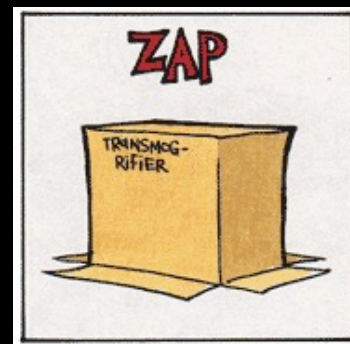
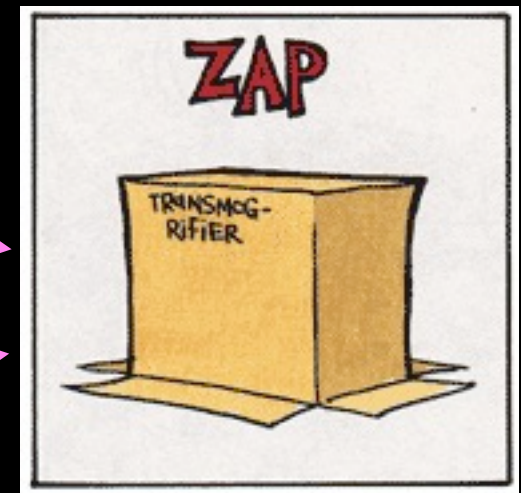
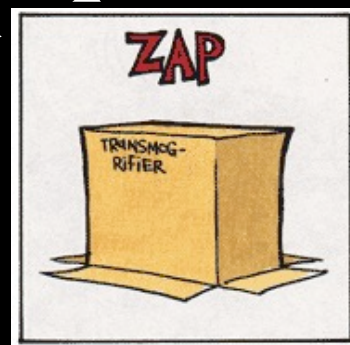
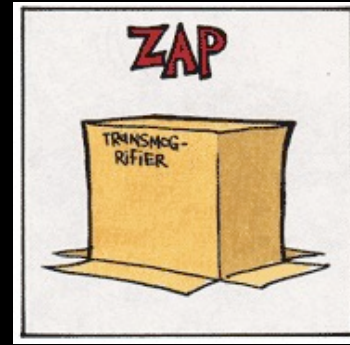
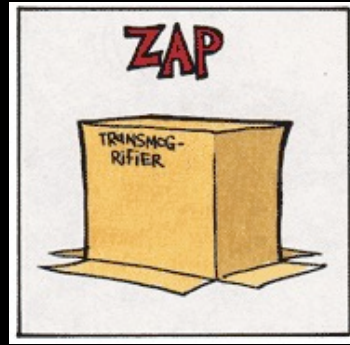
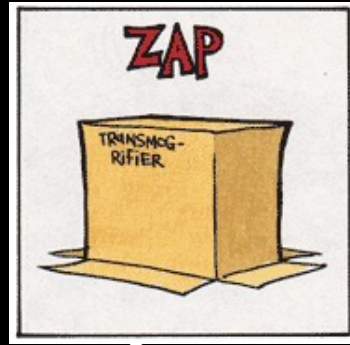
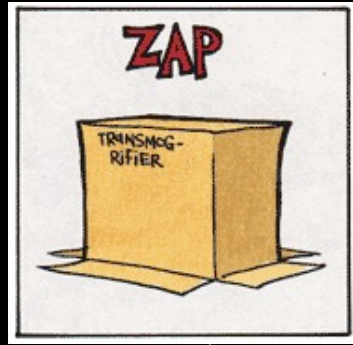
$$((4+16)+42)+(8+23)$$

12:00

12:01

12:02

12:03



12:00

12:01

12:02

12:03

```
8 module Stripe
9   def process_payment(pmt)
10     stats.time("transfer_ms") do
11       transfer(pmt.from, pmt.to, pmt.amount)
12     end
13   end
14 end
```

Maxifier




```
18 class Maxifier
19   def initialize
20     @result = 0
21   end
22
23   def update(x)
24     @result = [@result, x].max
25   end
26 end
27
```



12:00

4 8

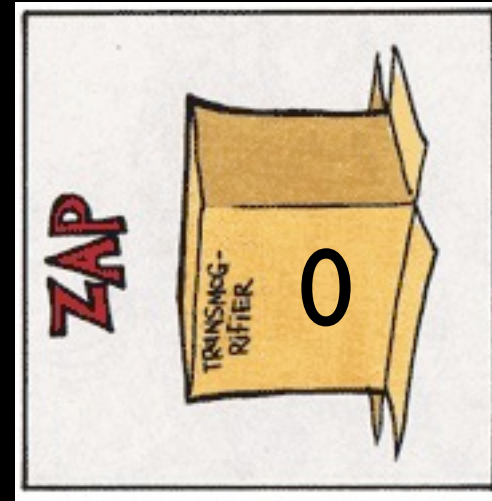


12:01

16 23



12:02



12:03

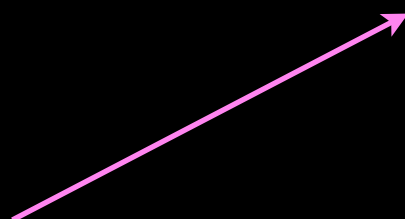
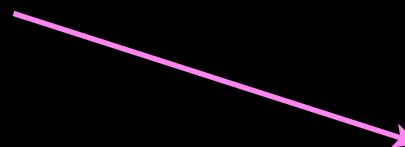
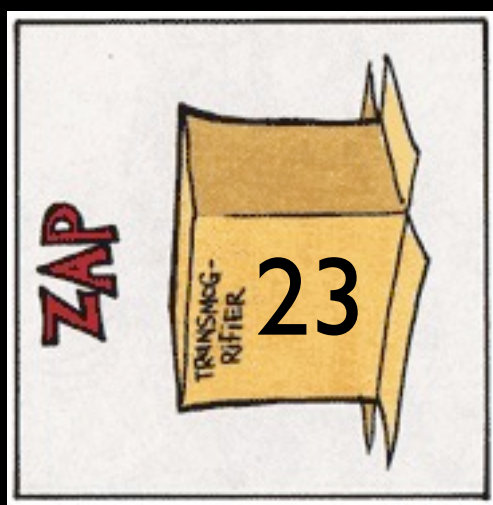
42



4 16 42



8 23

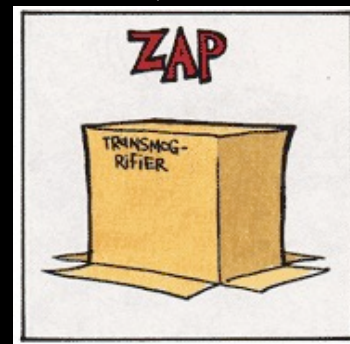
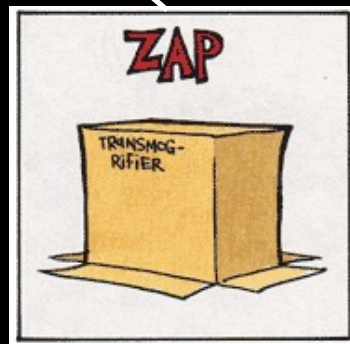
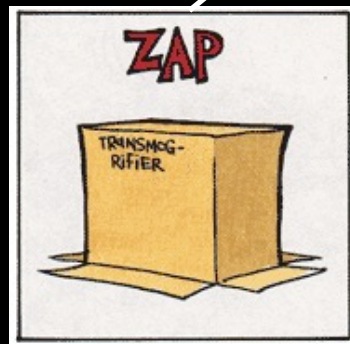
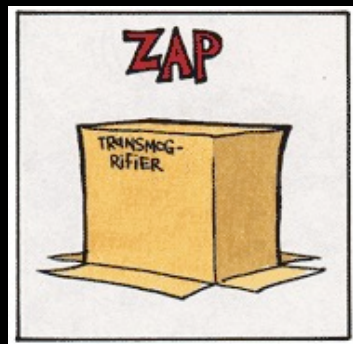
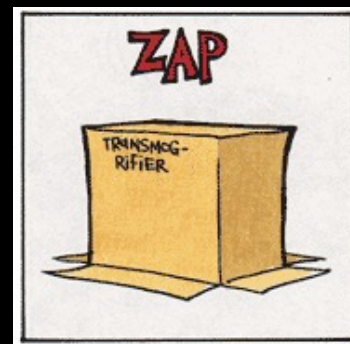
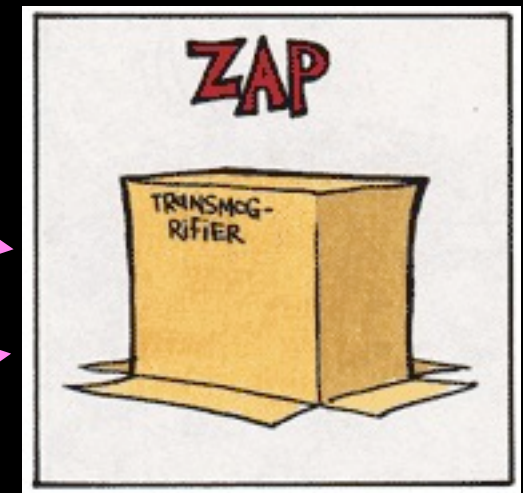
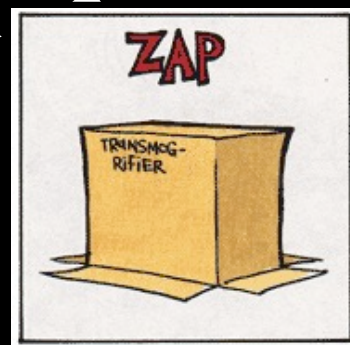
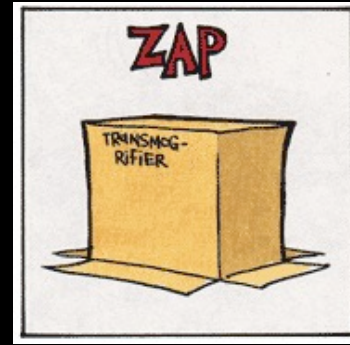
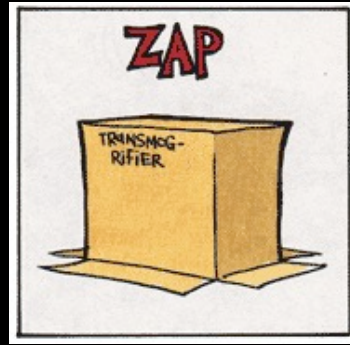
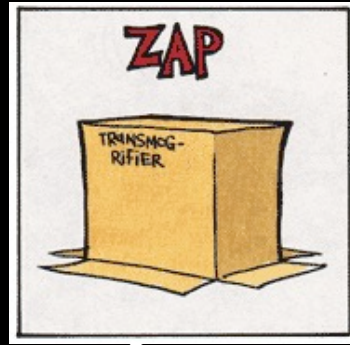
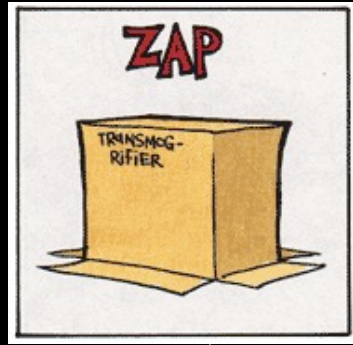


12:00

12:01

12:02

12:03



12:00

12:01

12:02

12:03

Generalizing + and max

- 1. Takes two numbers and produces another number
- 2. Grouping doesn't matter (associative)
- 3. Ordering doesn't matter (commutative)
- 4. Zeros get ignored

Commutative monoid

A set S , with an operation that:

- 1. Takes two members of S and produces a member of S
- 2. Grouping doesn't matter (associative)
- 3. Ordering doesn't matter (commutative)
- 4. Ignores some "identity" element of S

```

1 class TopK
2   def initialize(k)
3     @k = k
4     @result = {}
5   end
6
7   def update(hash)
8     merged = @result.merge(hash)
9     topk = merged.sort_by{|key,v| v}.reverse[0...@k]
10    Hash[topk]
11  end
12 end

```

{alice: 10}

{bob: 5}

{charlie: 7}

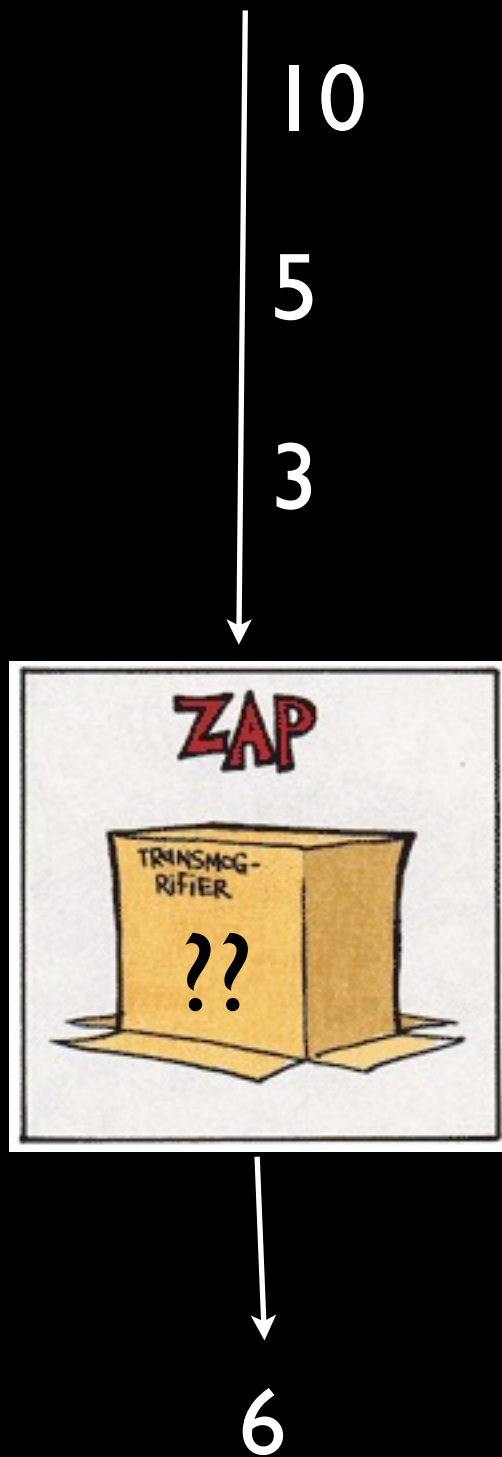
TopK Monoid



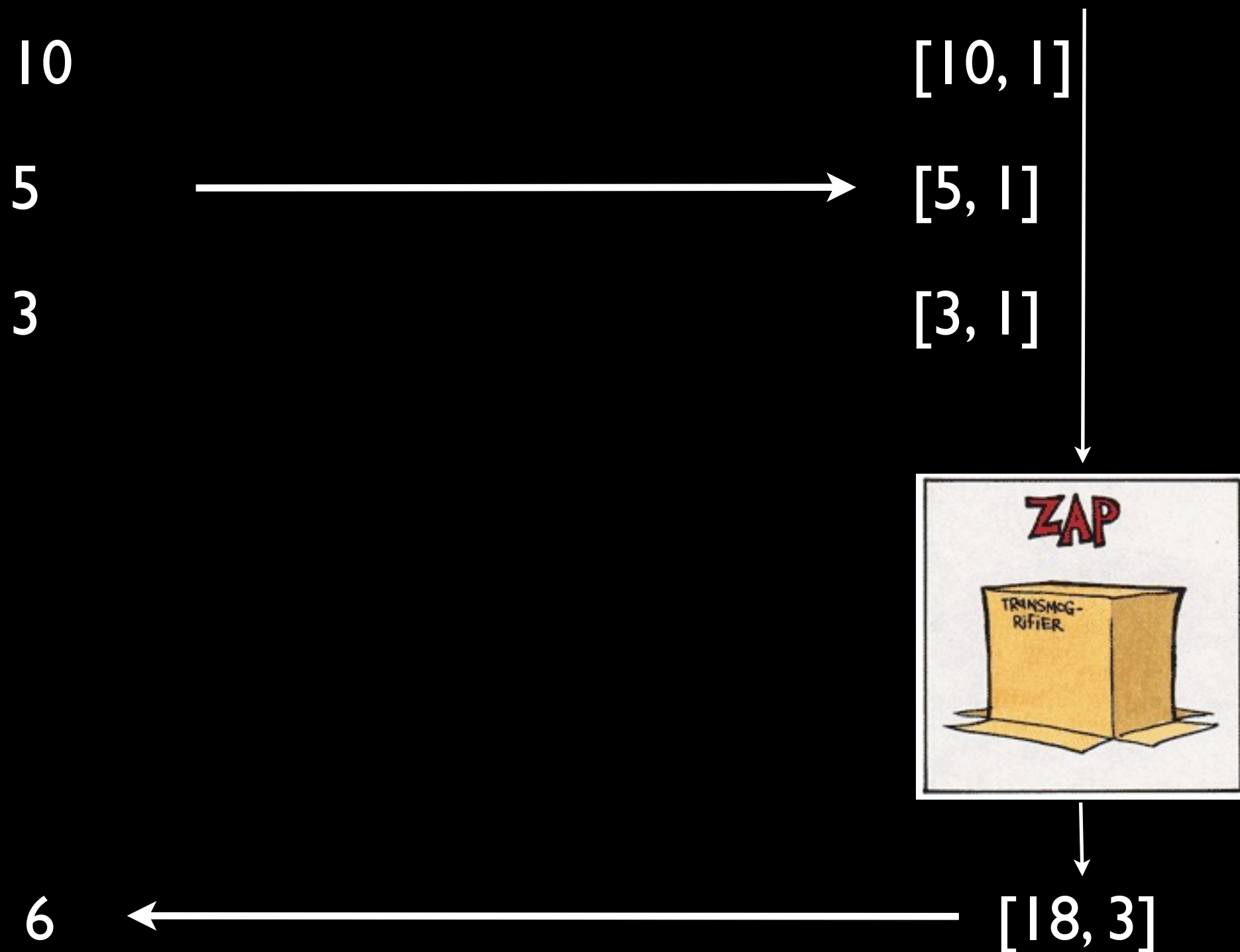
{alice: 10, charlie: 7}

(use a heap in real life, though)

Average Monoid



Average Monoid



(use a numerically stable average in real life, though)

Histogram Monoid

10

$[0,0,0,0,0,0,0,0,0,1]$

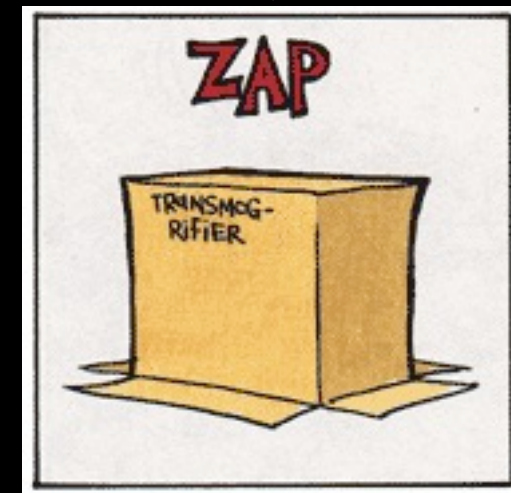
5



$[0,0,0,0,1,0,0,0,0,0]$

10

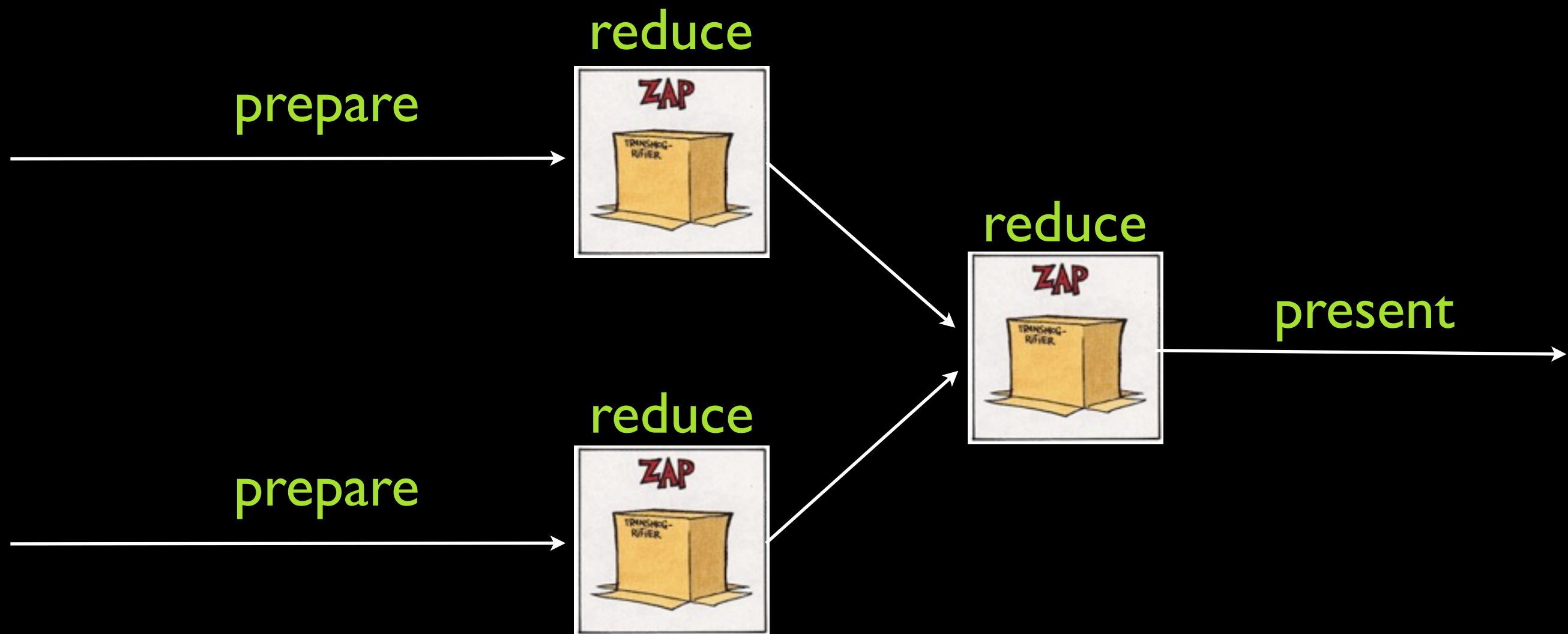
$[0,0,0,0,0,0,0,0,0,1]$



$[0,0,0,0,1,0,0,0,0,2]$



```
1 public interface Aggregator<I,0,S> {  
2     S initial();  
3     S prepare(I input);  
4     S reduce(S left, S right);  
5     0 present(S result);  
6 }  
7
```



Unique Values Monoid

alice

{alice}

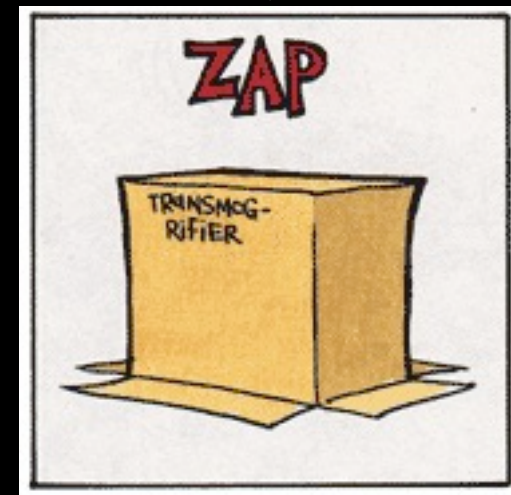
bob



{bob}

alice

{alice}



2



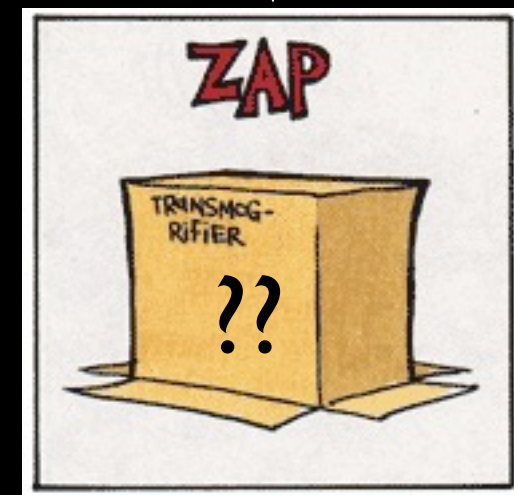
{alice,bob}

Unique Values Monoid

alice 0.789

bob $\xrightarrow{\text{hash}}$ 0.321

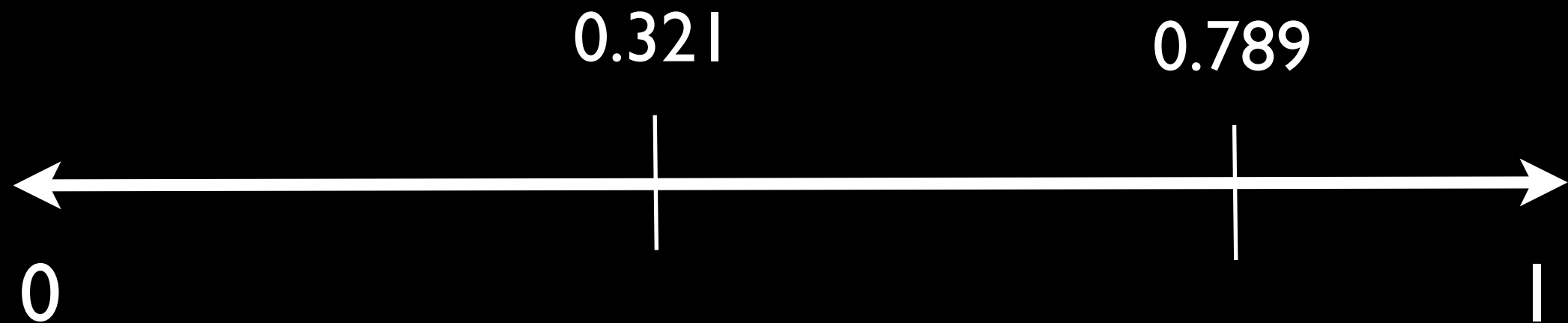
alice 0.789



2

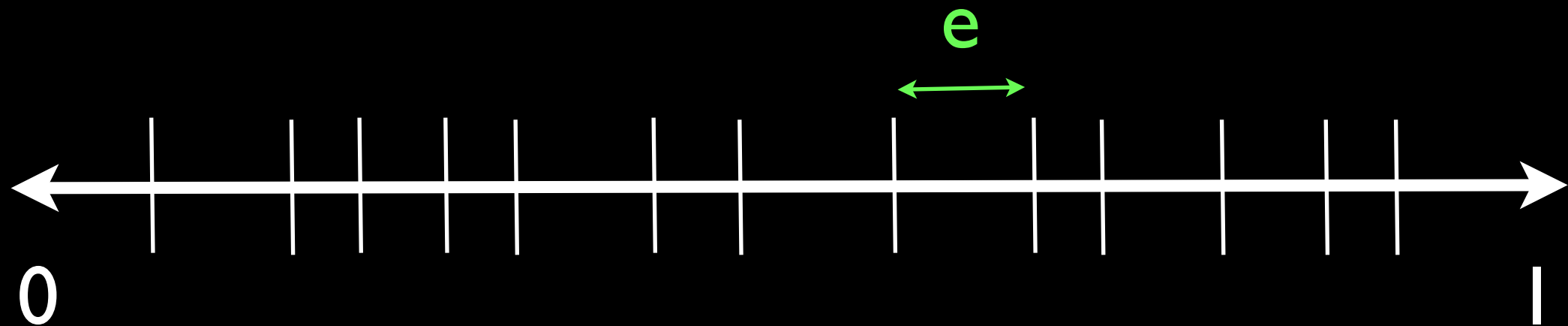


2 unique values



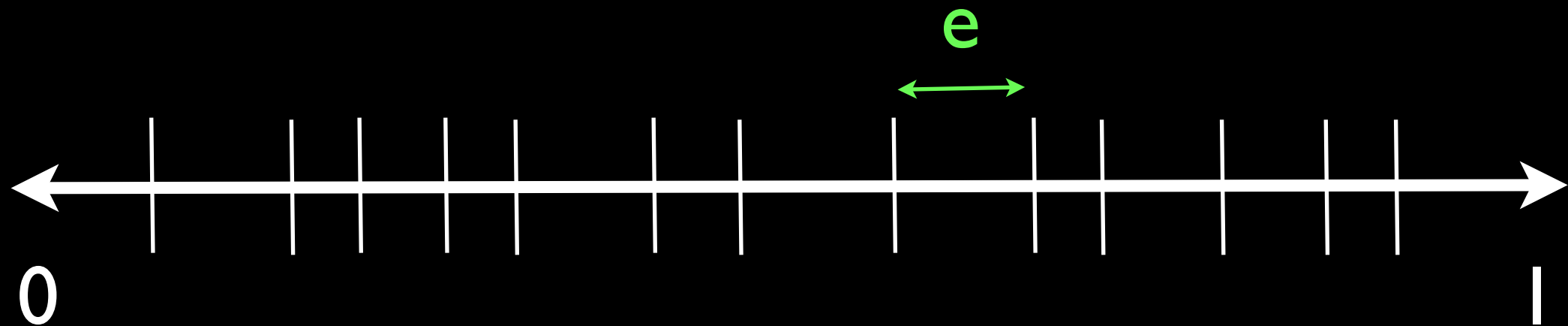
N unique values

$$E(e) = ??$$



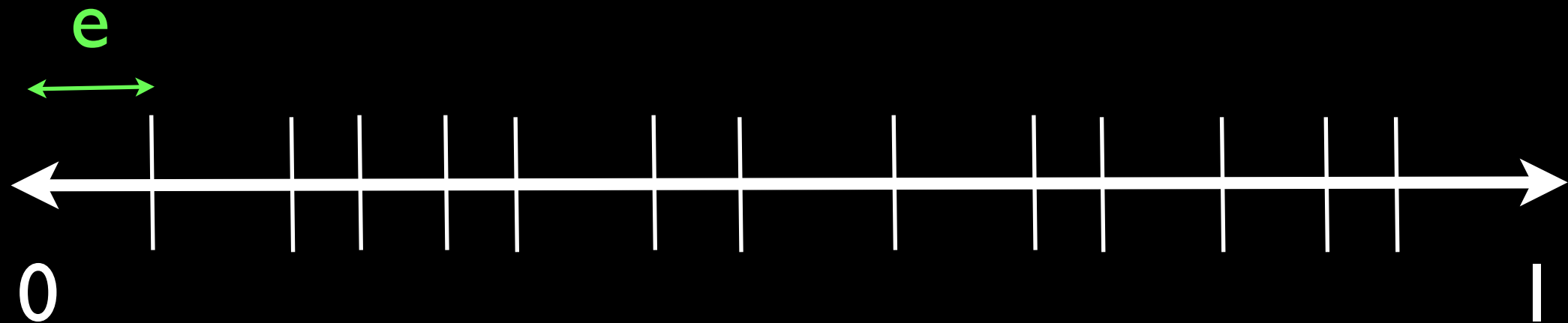
N unique values

$$E(e) = 1/(N+1)$$



N unique values

$$E(e) = 1/(N+1)$$



$$N = 1/e - 1$$

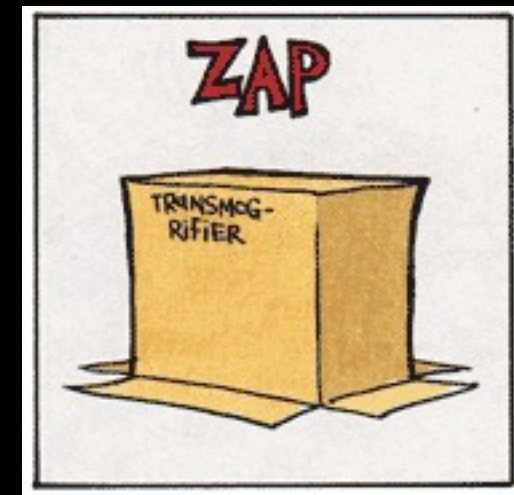
Unique Values Monoid

alice 0.789

bob $\xrightarrow{\text{hash}}$ 0.321

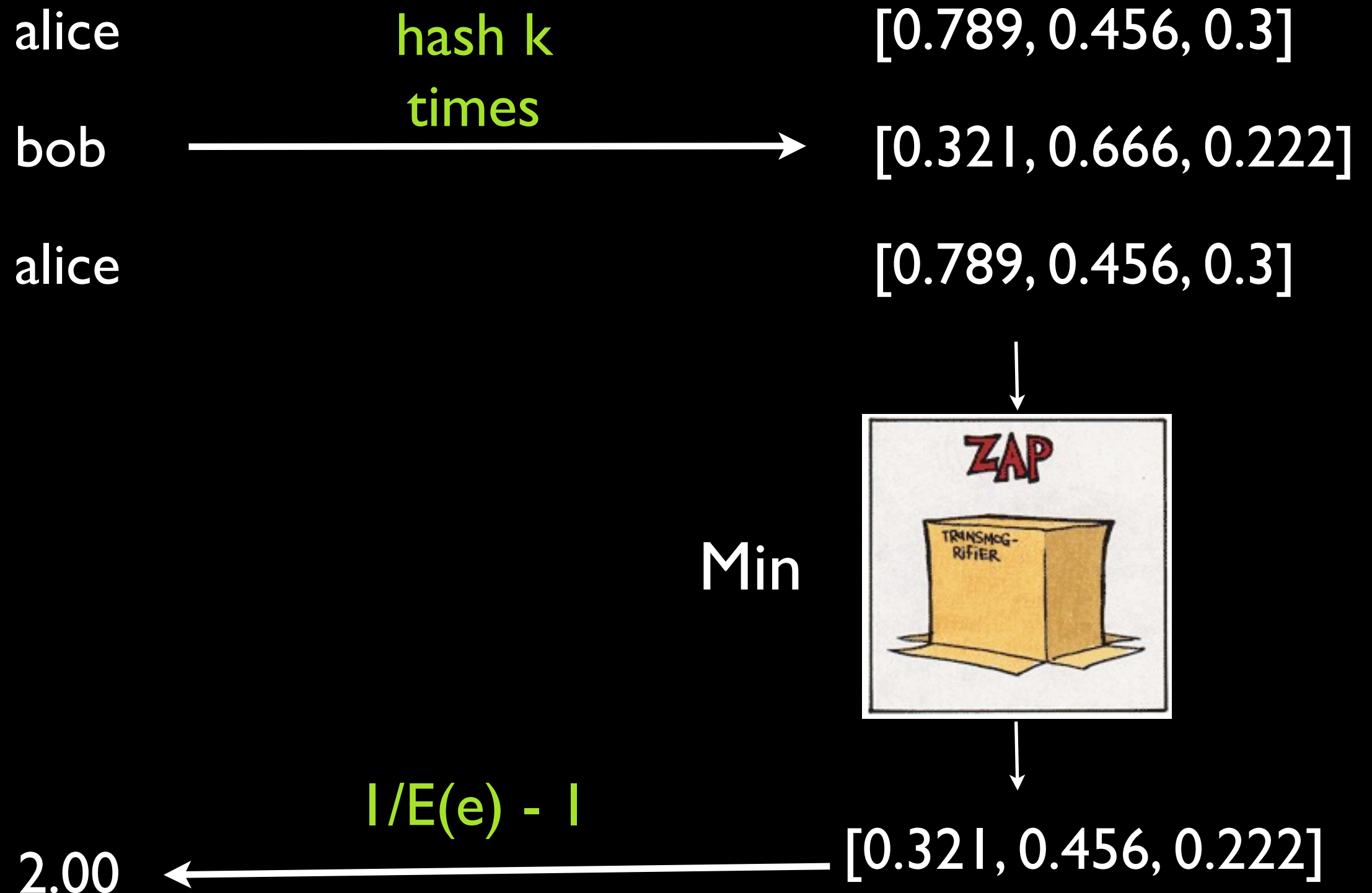
alice 0.789

Min



2.11 $\xleftarrow{|e| - 1}$ 0.321

Unique Values Monoid



In real life

- HyperLogLog for unique values
- Min-hash for set similarity
- Bloom filters for set membership

Frequency Monoid

alice

$\{\text{alice: } 1\}$

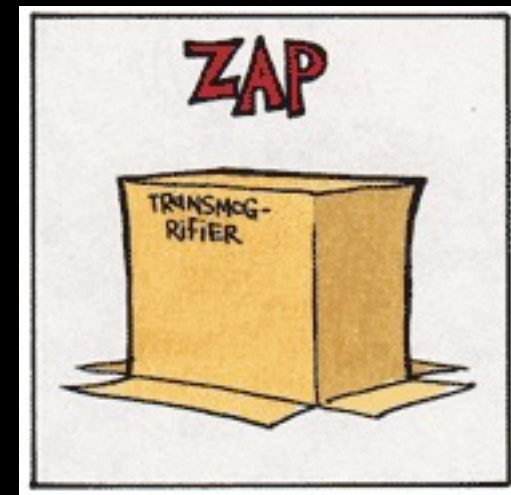
bob



$\{\text{bob: } 1\}$

alice

$\{\text{alice: } 1\}$



$\{\text{alice: } 2, \text{bob: } 1\}$

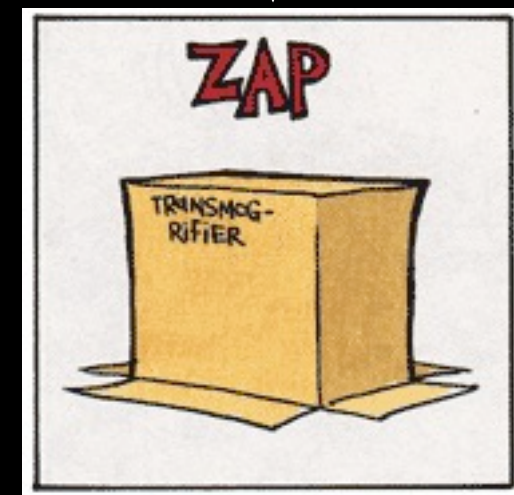
Frequency Monoid

alice $[0,0,1,0]$

hash % k

bob $[0,1,0,0]$

alice $[0,0,1,0]$



$[0,1,2,0]$

Frequency Monoid

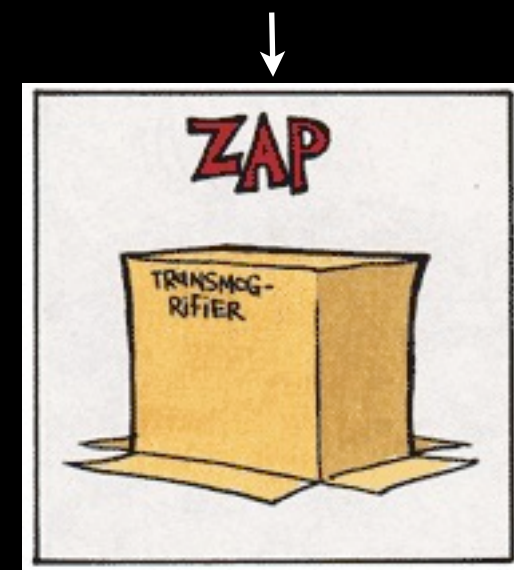
alice $[0,0,1,0]$

hash % k

bob $[0,1,0,0]$

alice $[0,0,1,0]$

charlie $[0,0,1,0]$



$[0,1,3,0]$

Count-min Sketch

alice

		2	
	2		
			2

Count-min Sketch

charlie

		3	
	2	1	
	1		2

Count-min Sketch

bob

	1	3	
	3	1	
1	1		2

Count-min Sketch

alice?

	1	3	
	3	1	
1	1		2

- Semigroup: set and associative operation
- Monoid: semigroup with identity
- Group: monoid with inverse

Any of these can be (and usually are) commutative

Commutative Monoids:

- Max
- HyperLogLog
- Bloom Filter
- ...

Abelian Groups:

- Sum
- Average
- Count-min Sketch
- ...

Subtraction!

- <http://github.com/twitter/algebird>
- <http://github.com/avibryant/simmer>
- <http://blog.aggregateknowledge.com>

@avibryant
avi@stripe.com