



iDigBio IT Standards Workshop
Gainesville, Florida

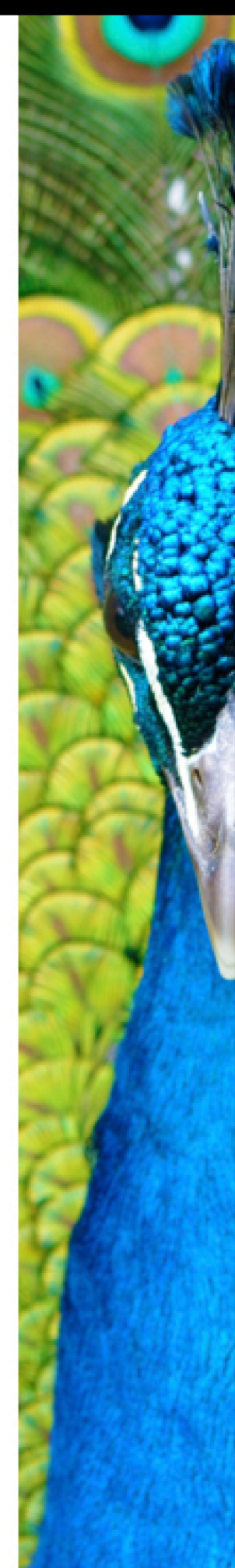
Scaling the Portal: A few lessons learnt

Tim Robertson
Systems Architect
Global Biodiversity Information Facility (GBIF)

27th March 2012

Standards: The *early* days

B i o d i v e r s i t y
I n f o r m a t i o n
S t a n d a r d s
T D W G

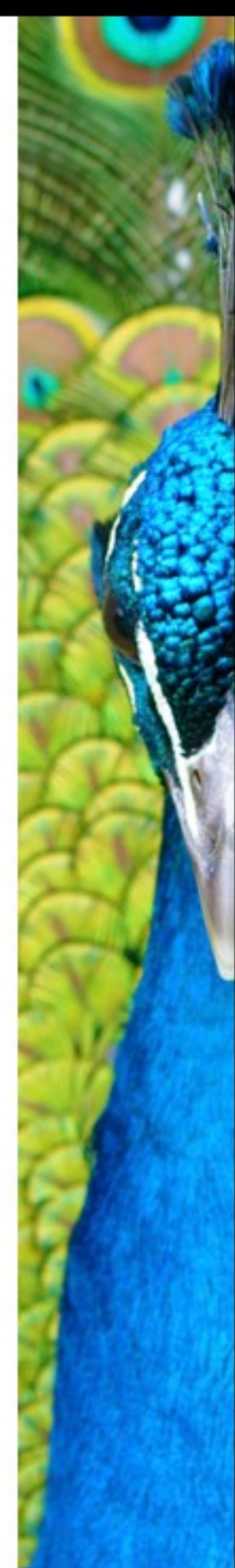


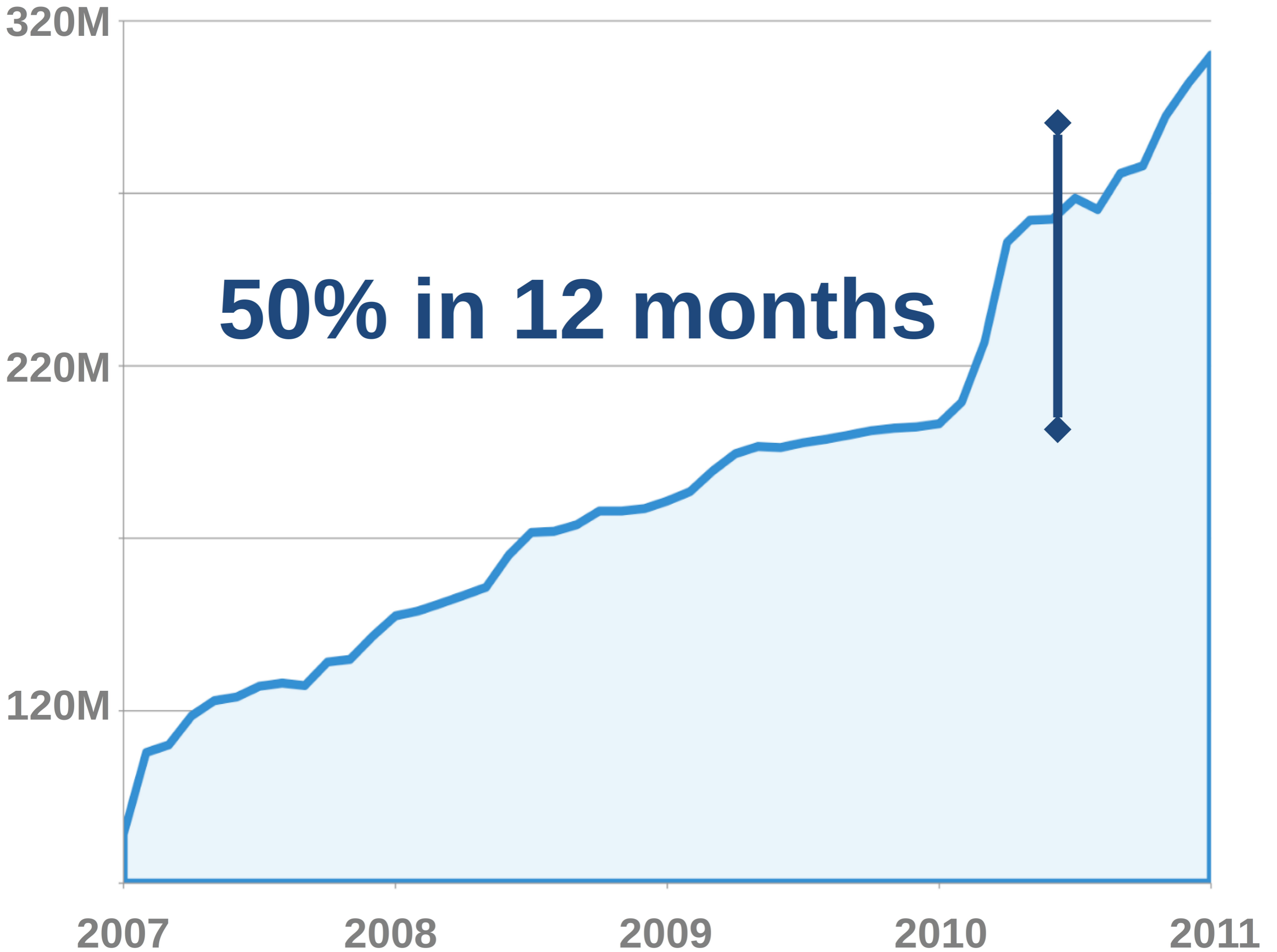
A pull model

- 1) Inventory of scientific names
- 2) Page by name range

Issues

- i. Verbose (latency)
- ii. `dwc:dateLastModified` (no deleted)
- iii. Full harvest needed...
- iv. Fragile
- v. Unique identifiers lacking





```
mysql>
```

```
SELECT SUM(c)
```

```
FROM
```

```
(SELECT COUNT(*) AS c FROM raw_occurrence_record UNION  
SELECT COUNT(*) AS c FROM occurrence_record UNION  
SELECT COUNT(*) AS c FROM identifier_record UNION  
SELECT COUNT(*) AS c FROM image_record UNION  
SELECT COUNT(*) AS c FROM typification_record UNION  
SELECT COUNT(*) AS c FROM link_record UNION  
SELECT COUNT(*) AS c FROM taxon_name UNION  
SELECT COUNT(*) AS c FROM taxon_concept UNION  
SELECT COUNT(*) AS c FROM common_name UNION  
SELECT COUNT(*) AS c FROM gmap_taxa_tile UNION  
SELECT COUNT(*) AS c FROM cell_density UNION  
SELECT COUNT(*) AS c FROM centi_cell_density) t1;
```

```
+-----+
```

```
| SUM(c) |
```

```
+-----+
```

```
| 2343864906 |
```

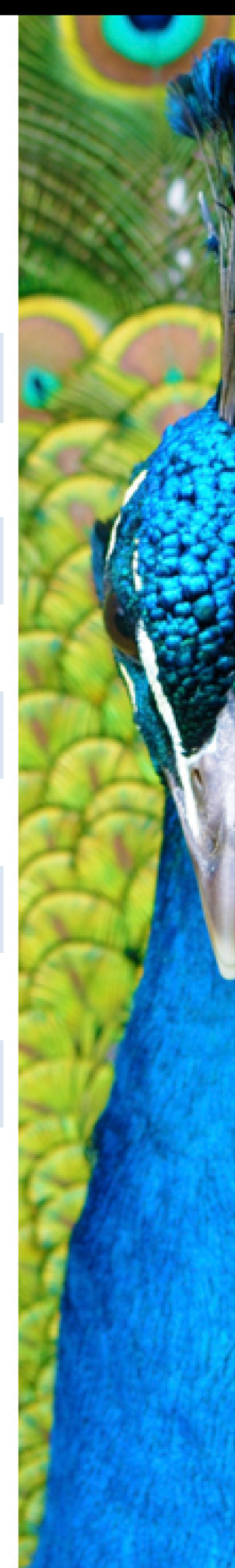
```
+-----+
```

```
1 row in set (0.01 sec)
```

Load on the data store

Table	Operations
raw_occurrence_record	2
taxon_name	14
taxon_concept	14
identifier_record	> 4
link_record	< 2
typification_record	< 2
image_record	0 - 10+
occurrence_record	1
Total	39+

**1 month processing
= 3000 Ops / sec.**



Harvester> "Delete records from dataset 123"

mysql> "Ok, I'm on it Sir!"

-- meanwhile

Harvester> "Perform an update of verbatim record"

Harvester> "Have you seen this record before?"

TaxaProc> "Extract distinct classifications"

GeoProc> "Extract distinct locations"

CountryProc> "Update country count"

PublisherProc> "Update publisher count"

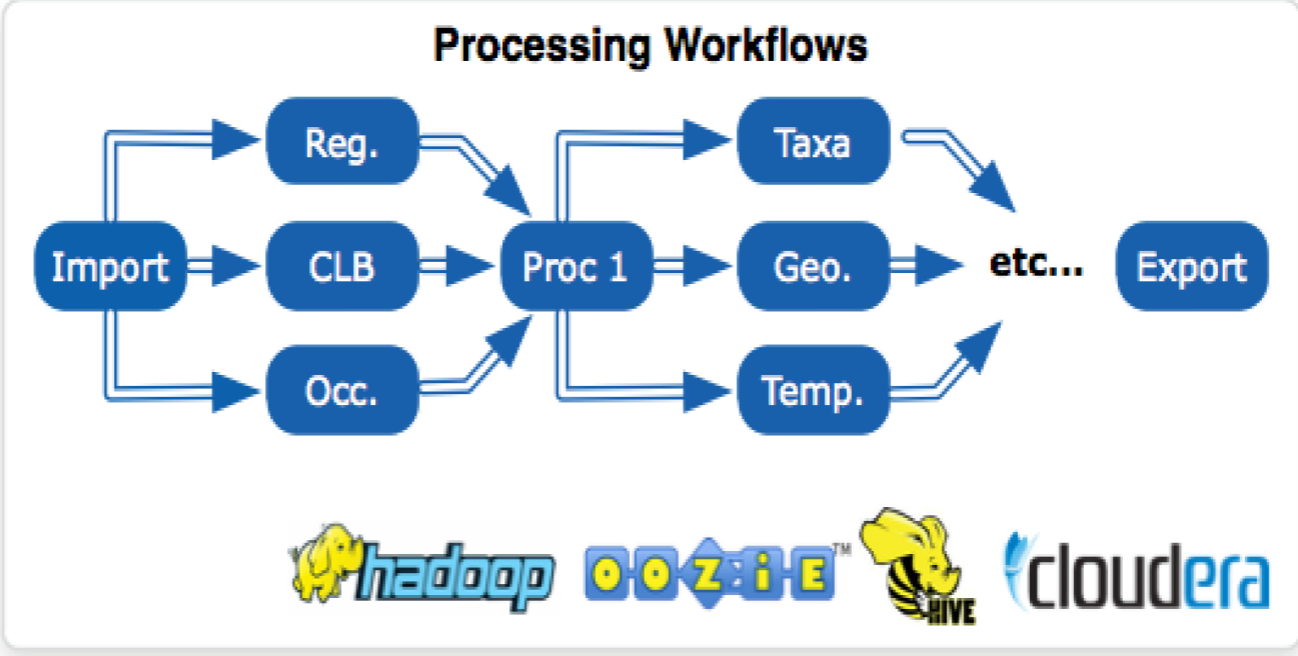
... etc etc

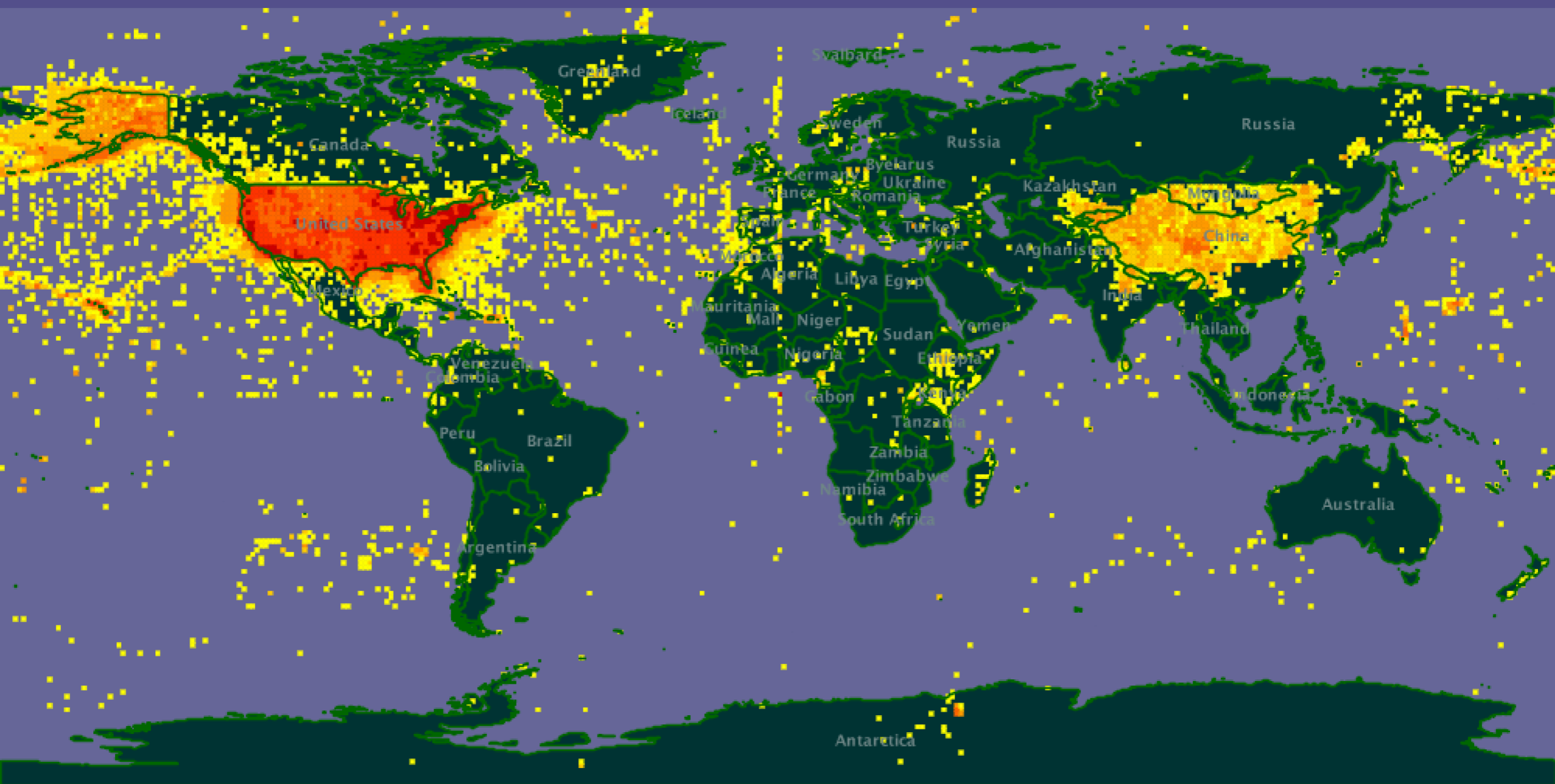
mysql> "Sorry guys, I'm locked right now"

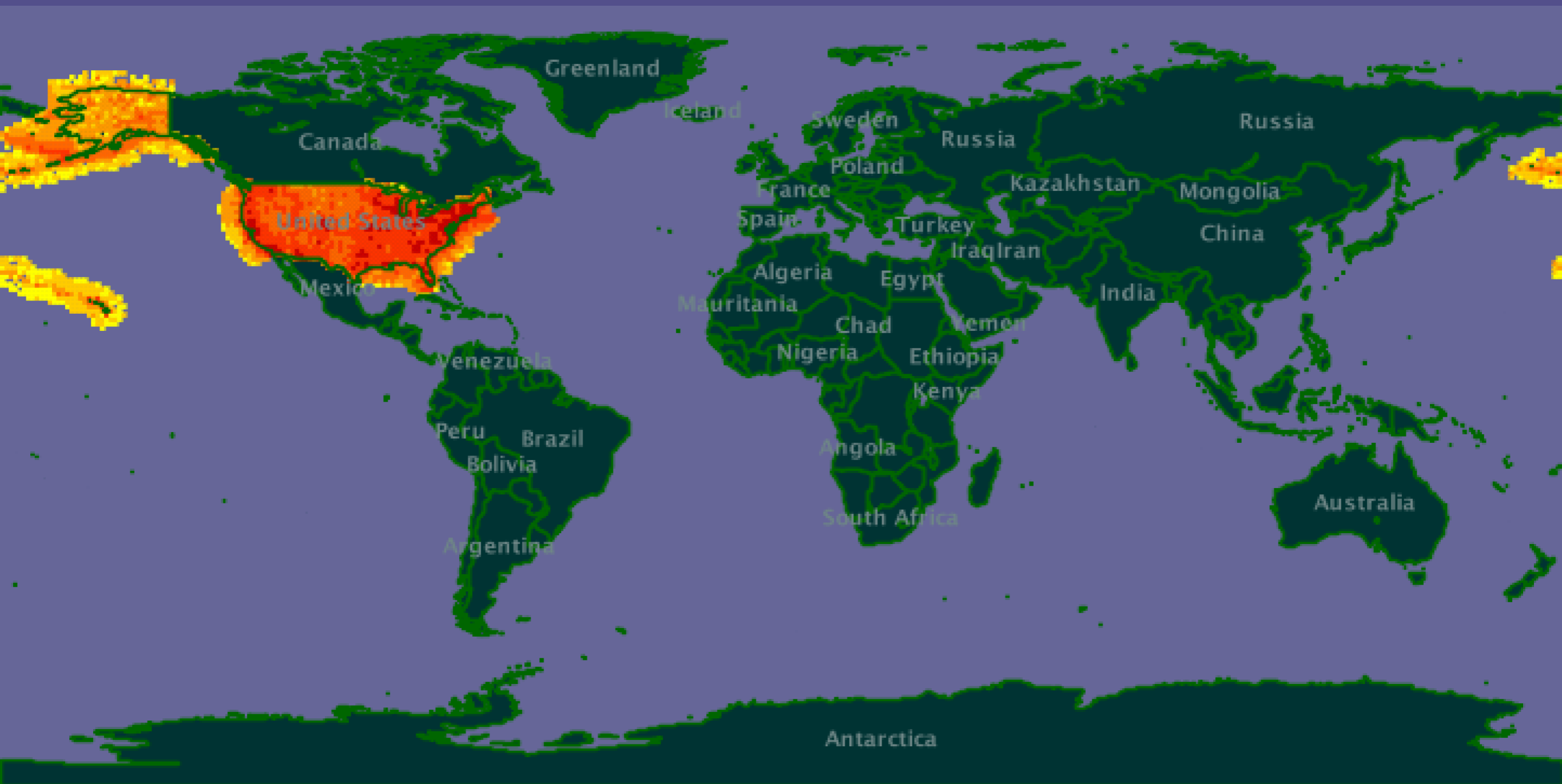


Harvesting and Indexing

Web applications and services







So what is the point?

- As we scale we need to reduce our work
 - i. Uniquely identify the record as soon as possible
 - ii. Detect if there has been a change
- Ability to track changes reduces latencies
 - i. Tracking will fail (message durability)
 - ii. Checkpoints are needed
 - iii. Transfer of full dataset
- Understand your data store technology
 - i. Locking can be a killer
 - ii. Understand query patterns
 - iii. Updating many indexes can be a killer
 - iv. What are the consistency requirements?

