

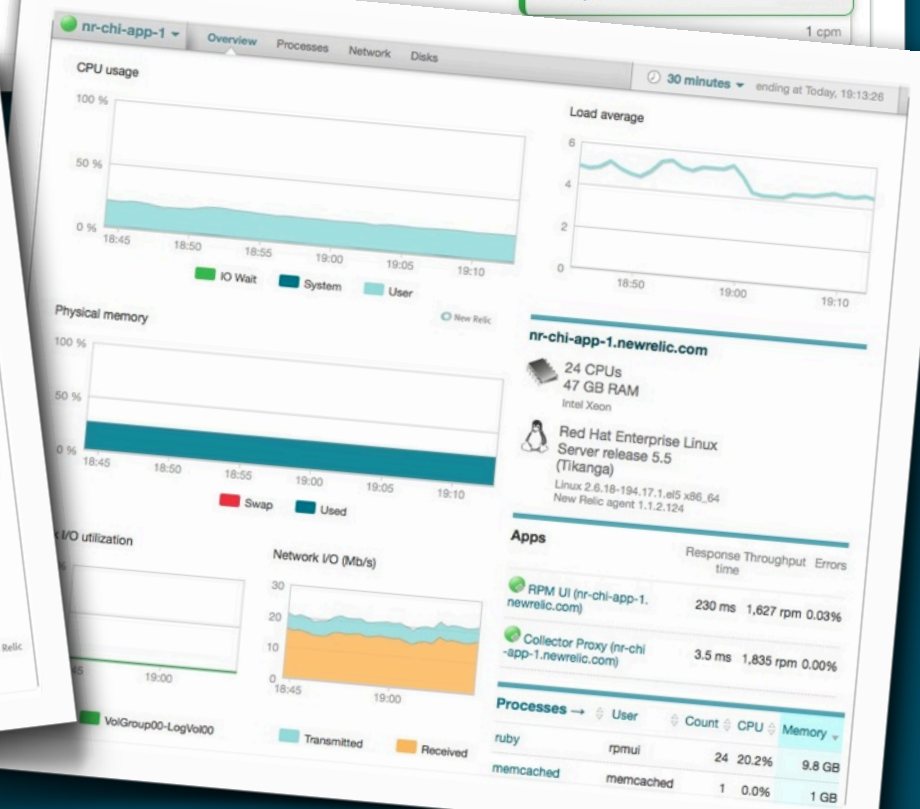
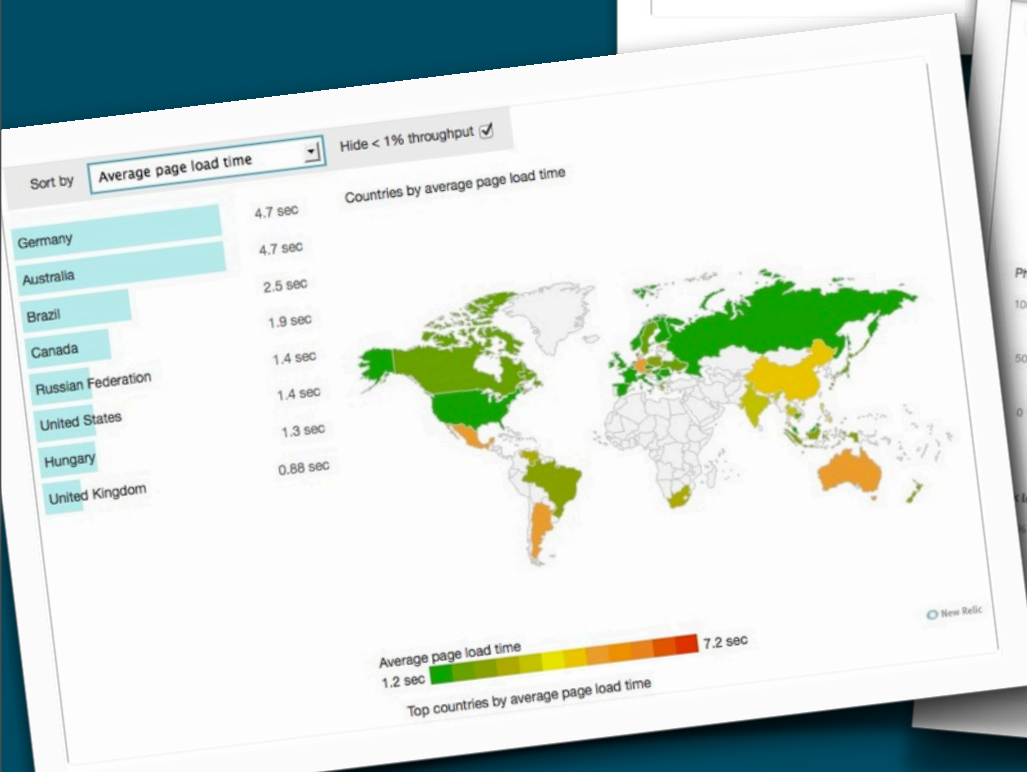
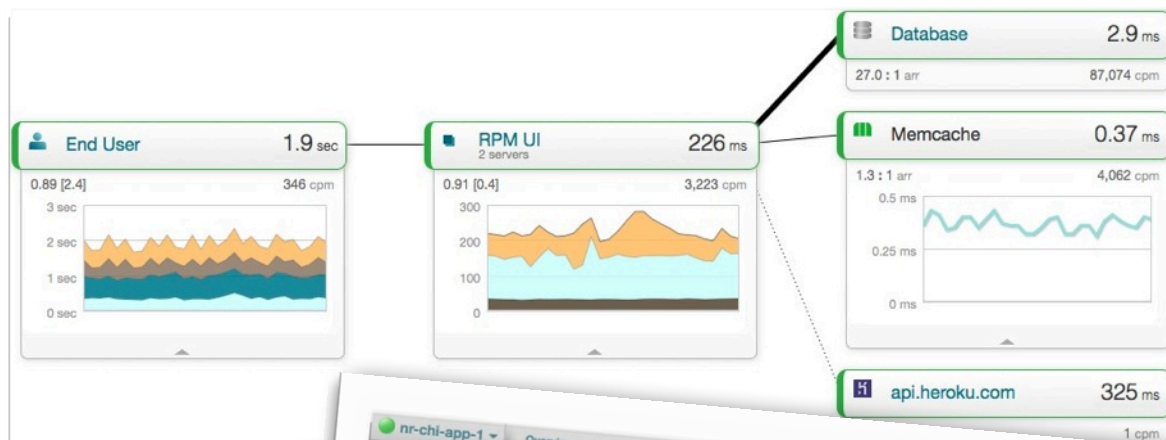
# SOME CONSIDERATIONS FOR SCALING

**Bjorn Freeman-Benson**  
*New Relic*

@bjorn\_fb



# New Relic

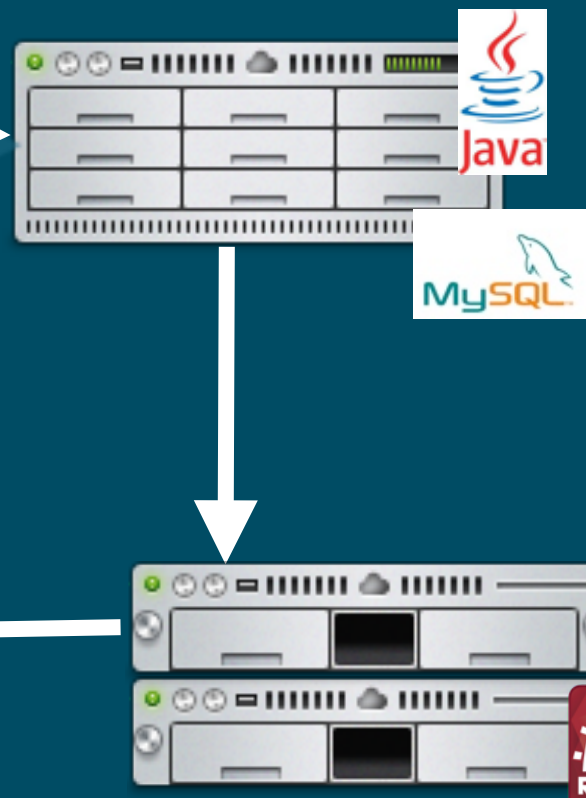




# Simplified Architecture

Our datacenter

Customer's environment



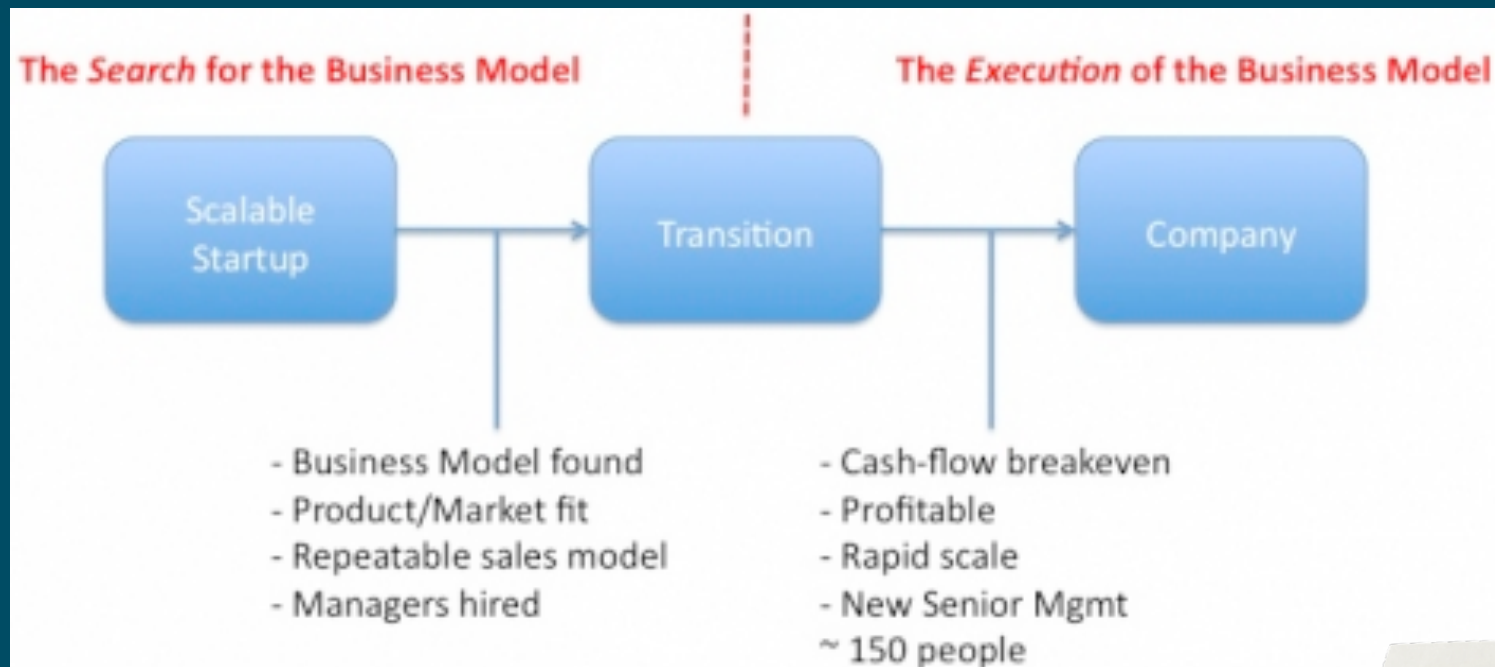
# Our growth

- In 4½ years, zero to 30,000 accounts...
- ... largest account has 17,000 servers
- ...  $58 \times 10^9$  metrics per day  
( $40 \times 10^6$  per minute)
- ... 5Tb of data a day  
(3.5Gb per minute)

[http://bit.ly/newrelic\\_stats](http://bit.ly/newrelic_stats)

# Lean Startup

- As a start-up: first prove that we had something, then scale, but plan to scale



<http://bit.ly/PPj3Yi>



# Our First System

- PaaS at Engine Yard
- 8 physical machines with multiple VMs
- Everything in Ruby
- Homegrown load balancer
- Separate processes for each activity
- **Perfect for the “Search for Business Model”**

# System Characteristics

- 1. Every app instance of every customer sends us data every minute**
- 2. Only a subset of customers view the data on any given minute**
- 3. Data has a steep half-life: most interesting data is seconds old**
- 4. Accuracy is essential**

# The Basics (5)





# #1: F5

- Reduce the number of connections to the servers
  - F5 buffers requests and handles SSL



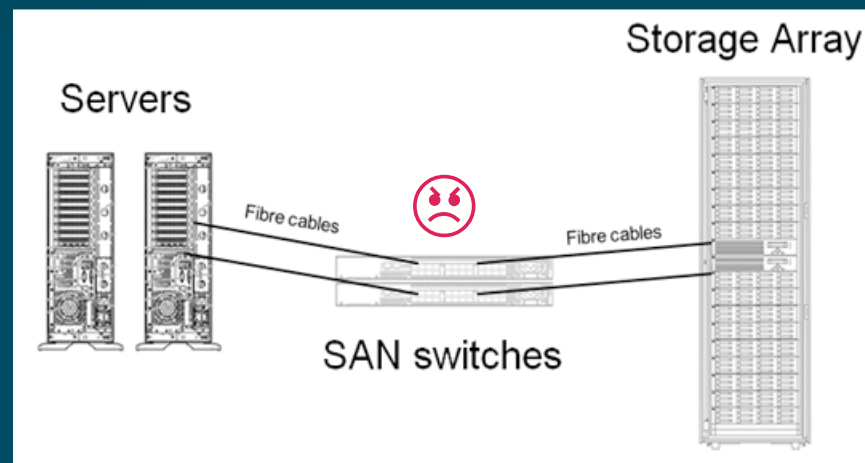
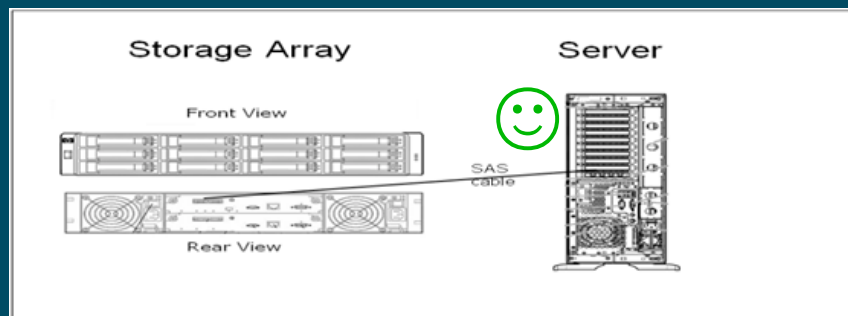
## #2: Bare metal

- VMs didn't work well for us
- I/O latency problems
- I/O bandwidth jitter
- Ruby is very memory heavy and VMs don't handle memory mapping as well as native CPUs

<http://bit.ly/Stmu5t>


# #3: Direct Attached Storage

- MySQL depends on really fast write commits
- Thus we need the disk cache as close to the cpu as possible



# #4: No App Servers

- Our high throughput collector processes don't need app servers so they are native Java apps with an embedded Jetty

 Aggregator		3.1 ms	598,771 rpm	0.01%	 ▼
 Beacon		0.17 ms	990,212 rpm	0.00%	 ▼

**jetty://**

<http://bit.ly/QrOExM>

# #5: Unicorn

- Every worker shares the socket so there's no need for a dispatcher
- Also easy to live-deploy new code - helps with our Continuous Deployment



<http://bit.ly/UAPARX>

# The Usual Suspects (4)





# #6: Agent Protocol

- Our first agent protocol was quick and dirty: Ruby object serialization and multiple round trips

```
def marshal_data(data)
  NewRelic::LanguageSupport.with_cautious_gc do
    Marshal.dump(data)
  end
rescue => e
  log.debug("#{e.class.name} : #{e.message} when marshalling #{object}")
  raise
end
```

- Refined: reduce round-trips (package more data into the payload); keep-alive

# #7: Accumulate & Resend



- If a service is temporarily unavailable, accumulate and retry

```
recover_from_communication_error:
```

```
/*
```

```
 * We were unable to contact the collectors, so we need to add all of this data to  
 * time unit's pending data.
```

```
*/
```

```
nr__log (NRL_DEBUG, "[%s] recovering from communication error..", appname);
```

```
nr__close_connection_to_daemon (nrdaemon);
```

```
nrthread_mutex_lock (&app->lock); {
```

```
/* merge metrics sets the ->replacement pointers of every metric in both from and
```

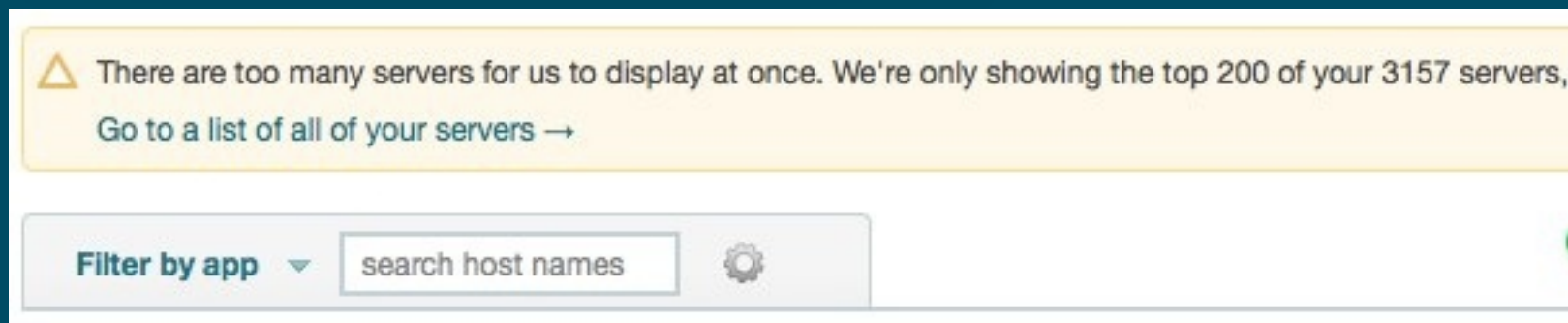
```
nr_metric_table__merge_metrics_from_to (data->metrics, app->pending_harvest->metr
```

```
nr__merge_slow_transactions_from_to (&(data->slow_transactions), &(app->pending_h
```

```
nr__merge_errors_from_to (&(data->errors), &(app->pending_harvest->errors));
```

# #8: Large Accounts

- Our first customers were small.
- Later larger customers stretched our assumptions. We added smart sorting, searching, paging, etc.



# #9: ORM Issues

- ORMs (Rails) are nice but can quickly load too many objects. Do a careful audit of slow code.

Slow transactions →	Resp. Time
ChartData::MetricChartsController#app_breakd... 12:16 — 30 minutes ago	435 ms
Api::V1::DataController#multi_app_data 12:16 — 30 minutes ago	2,230 ms
ApplicationsController#index 12:16 — 30 minutes ago	527 ms
Api::V1::DataController#multi_app_data 12:16 — 30 minutes ago	1,343 ms
ApplicationsController#index 12:16 — 30 minutes ago	1,272 ms
Show all slow transactions →	



# The Clever Stuff (6)

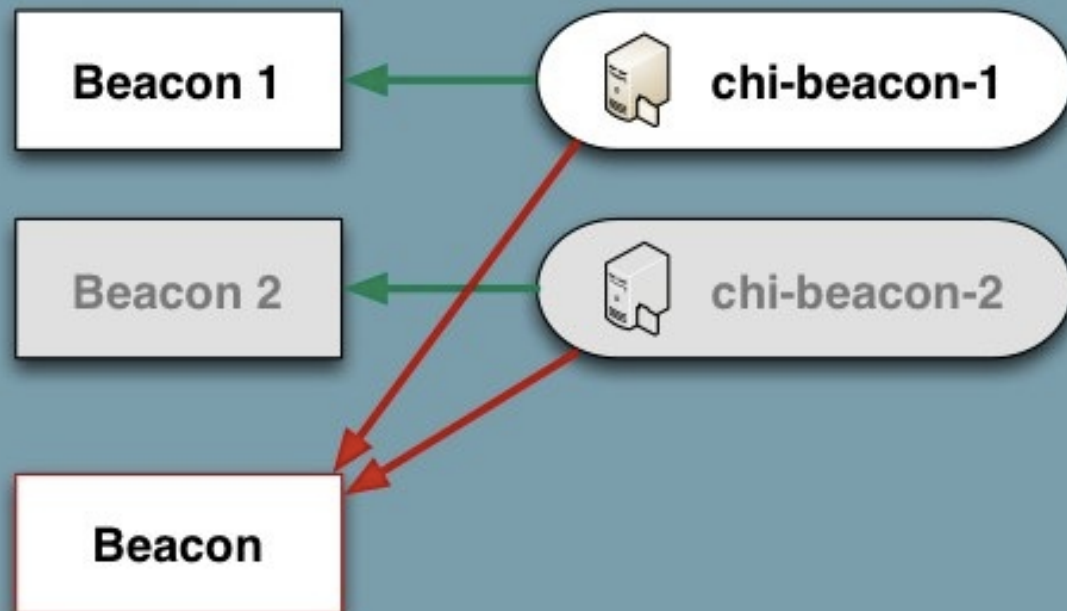


**UNITED FEDERATION OF AWESOMENESS**

# #10: Pre-compute

- Pre-compute expensive queries

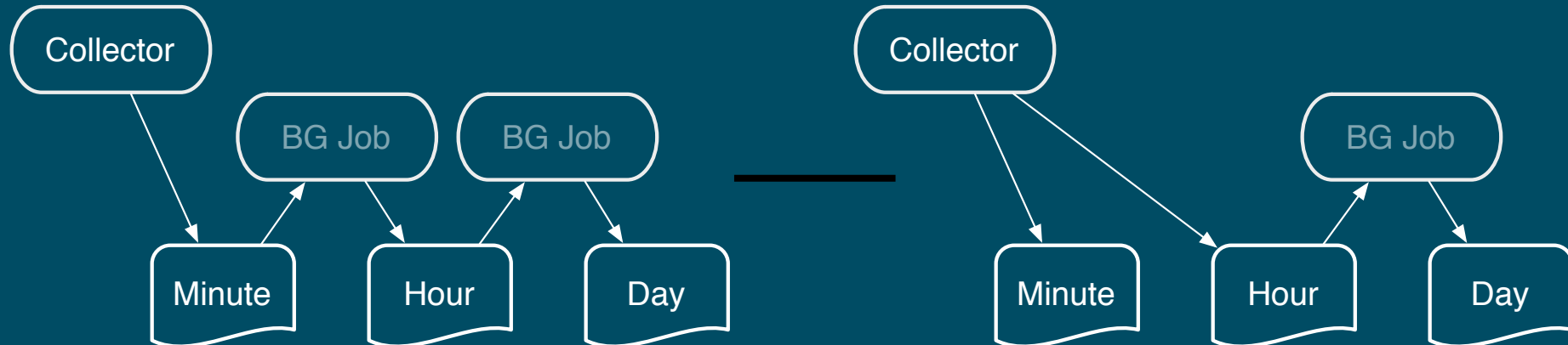
Beacon	0.2 ms	935,975 rpm	0.01%	⚙️
Beacon 1	0.22 ms	724,634 rpm ▼	0.01%	⚙️
Beacon 2	0.14 ms	211,331 rpm	0.00%	⚙️
Aggregator	8.8 ms			
Aggregator 0	1.9 ms			
Aggregator 1	5.3 ms			
Aggregator 2	3.8 ms			
Aggregator 3	4.1 ms			





# #11: Real-time BG

- **Background job to roll-up timeslice data: minutes to hours, hours to days**



# #12: Different DBs

- **Different data has different characteristics**
  - Account data is classic relational
  - Timeslice data is write-once
- **Use different database instances for each kind of data**
  - Different tuning parameters (buffer pools, etc)
  - **Similar to buddy memory allocation**

<http://bit.ly/VfQG8R>

# #13: Non-gc gc

- **Problem:** Deleting rows is expensive  
(due to table-level locking)
- **Solution:** Don't delete rows
  - Schema has multiple tables  
(one per account per time period)
  - Use DROP TABLE for gc
- Similar to the 100-request restart  
at [amazon.com/obidos](https://amazon.com/obidos) in 1999

# #14: Computation in DB

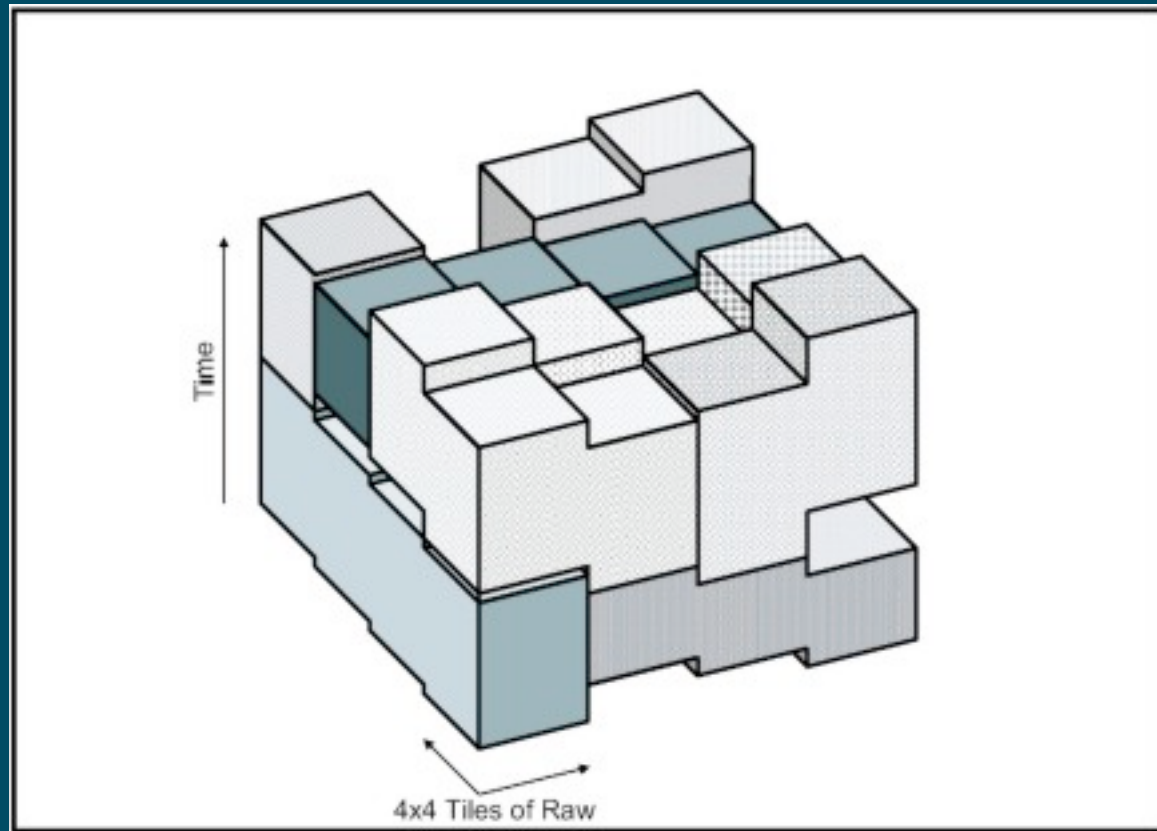
- **Natural sharding allows us to push computation into the db**
  - Supported by schema
  - Limits number of rows returned
  - Thus allows scripting language (Ruby) to do 'real' work
- **Opposite of the classical advice of doing nothing in the db**

<http://bit.ly/PFppZh>

# #15: SSDs

- **Our data is either random writes and sequential reads, or sequential writes and random reads**
  - Choose sequential reads because of UI
  - Use buffers to help random writes, but...
- **Switched to SSDs**
  - writes are same or slight slower
  - reads are fast, random or sequential

# The Optimizations (2)





# #16: Moving Processes

- Different processes have different performance characteristics: cpu, memory, i/o, time of day, etc.
- Allocate processes to machines to balance the resource requirements
  - Instead of “all X processes on M1 and Ys on M2” we balance the machines

# #17: Moving Customers

- **Customers have different data characteristics: size, access patterns, ...**
- **Allocate customers to shards to balance the size and loads on the shards**
  - **Required an early architectural decision to allow data split between shards**

# Take-away



Lessons  
Learned

# Take-away

1. Do the basics
2. Design in some scalability
3. Use the unique characteristics of your app to optimize
4. Buzzwords not needed

