# Facebook Architecture

Aditya Agarwal
Director of Engineering
11/22/2008

# Agenda

| | |
|---|---|
| **1** | Architecture Overview |
| **2** | PHP, MySQL, Memcache |
| **3** | Thrift, Scribe, Tools |
| **4** | News Feed Architecture |

# facebook At a Glance



## The Social Graph

120M+ active users

50B+ PVs per month

10B+ Photos

1B+ connections

50K+ Platform Apps

400K+ App Developers

# General Design Principles

- Use open source where possible
    - Explore making optimizations where needed
- Unix Philosophy
    - Keep individual components simple yet performant
    - Combine as necessary
    - Concentrate on clean interface points
- Build everything for scale
- Try to minimize failure points
- Simplicity, Simplicity, Simplicity!

# Architecture Overview

## LAMP    +    Services

PHP

memcached

MySQL

| PHP |
| --- |
| Memcache |
| MySQL |

php!

| AdServer |
| --- |
| Search |
| Network Selector |
| News Feed |
| Blogfeeds |
| CSSParser |
| Mobile |
| ShareScraper |

!php

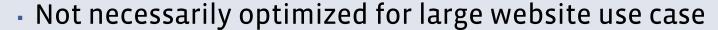| Thrift |
| --- |
| Scribe |
| ODS |
| Tools |

# PHP



- Good web programming language

    - Extensive library support for web development

    - Active developer community


- Good for rapid iteration

    - Dynamically typed, interpreted scripting language

# PHP: What we Learnt



- Tough to scale for large code bases
  - Weak typing
  - Limited opportunities for static analysis, code optimizations

- Not necessarily optimized for large website use case
  - E.g. No dynamic reloading of files on web server

- Linearly increasing cost per included file

- Extension framework is difficult to use

# PHP: Customizations



- Op-code optimization

- APC improvements
  - Lazy loading
  - Cache priming
  - More efficient locking semantics for variable cache data
- Custom extensions
  - Memcache client extension
  - Serialization format
  - Logging, Stats collection, Monitoring
  - Asynchronous event-handling mechanism

# MySQL

- Fast, reliable

- Used primarily as ‹key,value› store
    - Data randomly distributed amongst large set of logical instances
    - Most data access based on global id

- Large number of logical instances spread out across physical nodes
    - Load balancing at physical node level

- No read replication

# MySQL: What We Learnt (ing)

- Logical migration of data is *very* difficult

- Create a large number of logical dbs, load balance them over varying number of physical nodes

- No joins in production
  - Logically difficult (because data is distributed randomly)

- Easier to scale CPU on web tier

# MySQL: What we Learnt (ing)

- Most data access is for recent data
    - Optimize table layout for recency
    - Archive older data


- Don't ever store non-static data in a central db
    - CDB makes it easier to perform certain aggregated queries
    - Will not scale


- Use services or memcache for global queries
    - E.g.: What are the most popular groups in my network

# MySQL: Customizations

- No extensive native MySQL modifications


- Custom partitioning scheme
    - Global id assigned to all data


- Custom archiving scheme
    - Based on frequency and recency of data on a per-user basis


- Extended Query Engine for cross-data center replication, cache consistency

# MySQL: Customizations

- Graph based data-access libraries

  - Loosely typed objects (nodes) with limited datatypes (int, varchar, text)

  - Replicated connections (edges)

  - Analogous to distributed foreign keys

- Some data collocated

  - Example: User profile data and all of user's connections

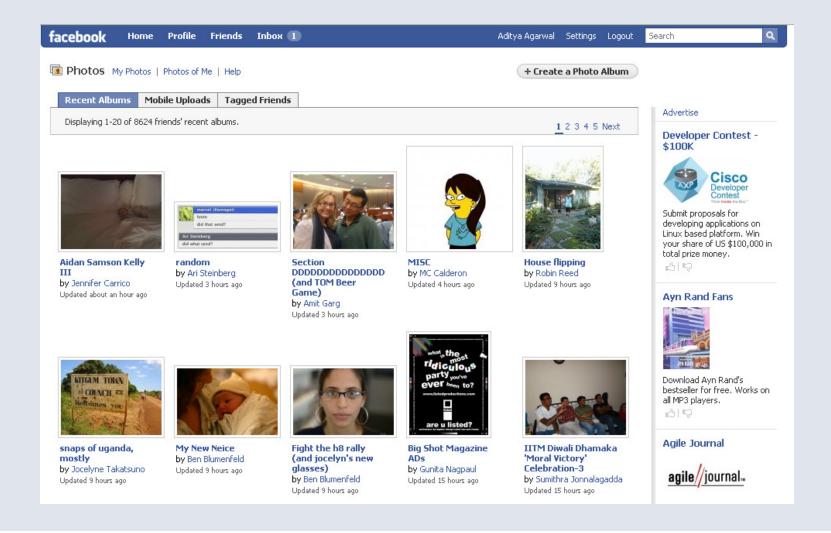- Most data distributed randomly

# Memcache

- High-Performance, distributed in-memory hash table

- Used to alleviate database load

- Primary form of caching

- Over 25TB of in-memory cache

- Average latency < 200 micro-seconds

- Cache serialized PHP data structures

- Lots and lots of multi-gets to retrieve data spanning across graph edges

# Memache: Customizations

- Memache over UDP
    - Reduce memory overhead of thousands of TCP connection buffers
    - Application-level flow control (optimization for multi-gets)

- On demand aggregation of per-thread stats
    - Reduces global lock contention

- Multiple Kernel changes to optimize for Memcache usage
    - Distributing network interrupt handling over multiple cores
    - Opportunistic polling of network interface

# Let's put this into action

# Under the Covers

- Get my profile data

  - Fetch from cache, potentially go to my DB (based on user-id)

- Get friend connections

  - Cache, if not DB (based on user-id)

- In parallel, fetch last 10 photo album ids for each of my friends

  - Multi-get; individual cache misses fetches data from db (based on photo-album id)

- Fetch data for most recent photo albums in parallel

- Execute page-specific rendering logic in PHP

- Return data, make user happy

# LAMP is not Perfect

# LAMP is not Perfect

- PHP+MySQL+Memcache works for a large class of problems but not for everything
    - PHP is stateless
    - PHP not the fastest executing language
    - All data is remote
- Reasons why services are written
    - Store code closer to data
    - Compiled environment is more efficient
    - Certain functionality only present in other languages

# Services Philosophy

- Create a service iff required

  - Real overhead for deployment, maintenance, separate code-base

  - Another failure point

- Create a common framework and toolset that will allow for easier creation of services

  - Thrift

  - Scribe

  - ODS, Alerting service, Monitoring service

- Use the right language, library and tool for the task

# Thrift









High-Level Goal: Enable transparent interaction between these.
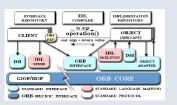
...and some others too.

# Thrift

- Lightweight software framework for cross-language development

- Provide IDL, statically generate code

- Supported bindings: C++, PHP, Python, Java, Ruby, Erlang, Perl, Haskell etc.

- Transports: Simple Interface to I/O

  - Tsocket, TFileTransport, TMemoryBuffer

- Protocols: Serialization Format

  - TBinaryProtocol, TJSONProtocol

- Servers

  - Non-Blocking, Async, Single Threaded, Multi-threaded

# Hasn't this been done before? (yes.)

- SOAP

  - XML, XML, and more XML

- CORBA

  - Bloated? Remote bindings?

- COM

  - Face-Win32ClientSoftware.dll-Book

- Pillar

  - Slick! But no versioning/abstraction.

- Protocol Buffers

# Thrift: Why?

- It's quick. Really quick.

- Less time wasted by individual developers
    - No duplicated networking and protocol code
    - Less time dealing with boilerplate stuff
    - Write your client and server in about 5 minutes

- Division of labor
    - Work on high-performance servers separate from applications

- Common toolkit
    - Fosters code reuse and shared tools

# Scribe

- Scalable distributed logging framework

- Useful for logging a wide array of data

    - Search Redologs

    - Powers news feed publishing

    - A/B testing data

- Weak Reliability

    - More reliable than traditional logging but not suitable for database transactions.

- Simple data model

- Built on top of Thrift

# Other Tools

- SMC (Service Management Console)

  - Centralized configuration

  - Used to determine logical service -> physical node mapping
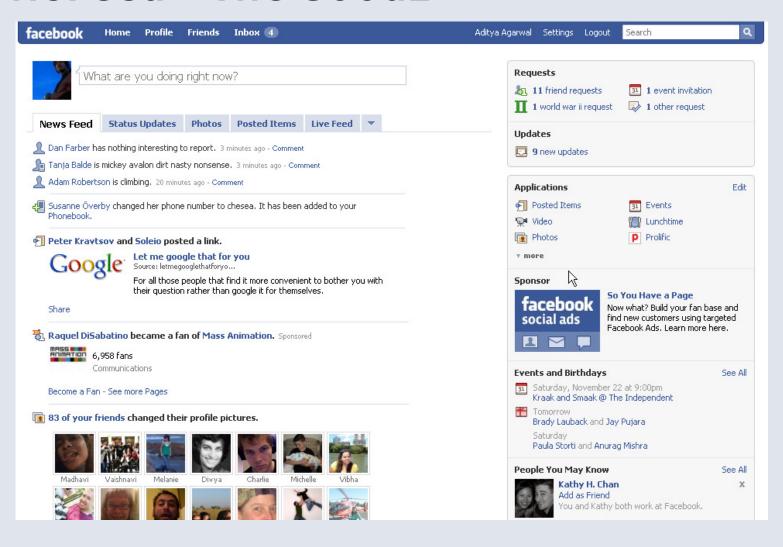
# Other Tools

- ODS

  - Used to log and view historical trends for any stats published by service

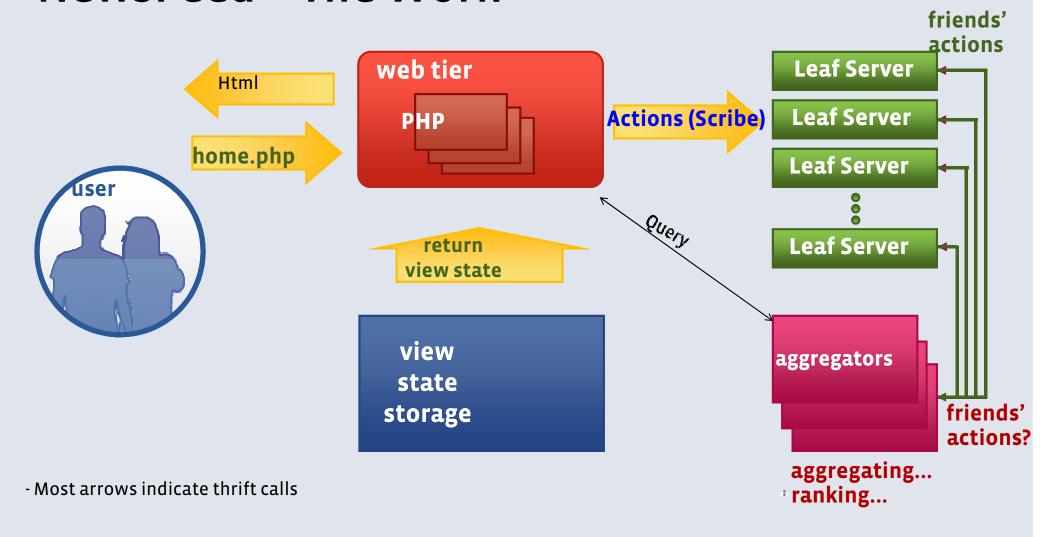  - Useful for service monitoring, alerting

# Open Source

- Thrift

  - http://developers.facebook.com/thrift/

- Scribe

  - http://developers.facebook.com/scribe/

- PHPEmbed

  - http://developers.facebook.com/phpembed/

- More good stuff

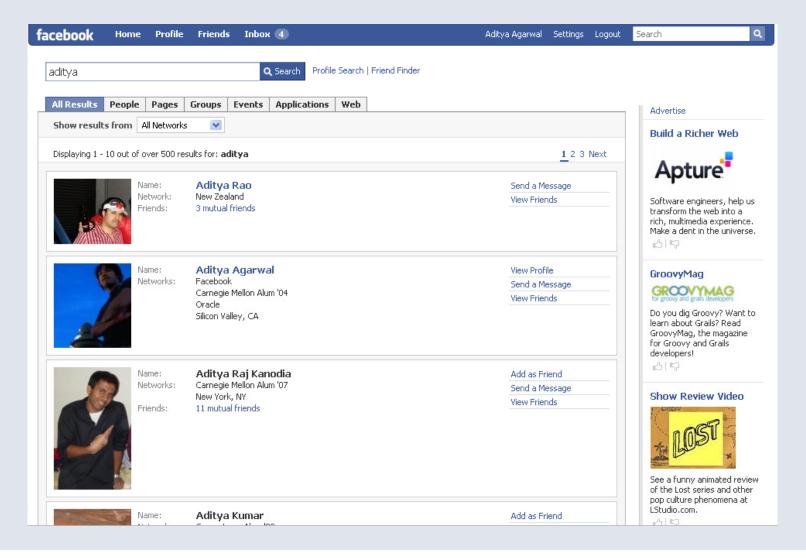  - http://developers.facebook.com/opensource.php
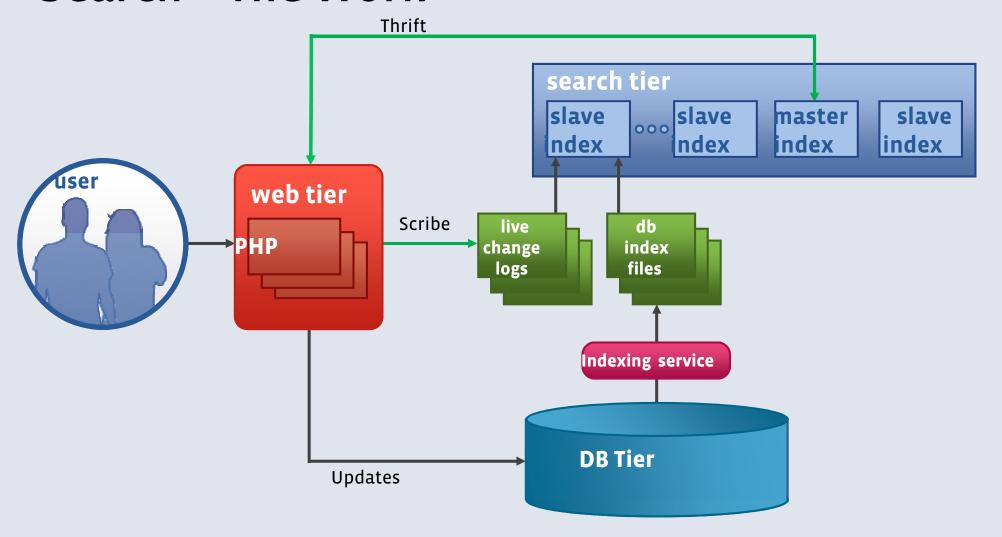
# NewsFeed – The Goodz

# NewsFeed – The Work

**user**

Html ←

home.php →

**web tier**

**PHP**

**Actions (Scribe)** →

**Leaf Server**

**Leaf Server**

**Leaf Server**

**Leaf Server**

friends'
actions

return
view state

Query

**view
state
storage**

**aggregators**

friends'
actions?

aggregating...
ranking...

- Most arrows indicate thrift calls

# Search – The Goodz

# Search – The Work

# Questions?

More info at www.facebook.com/eblog

Aditya Agarwal
aditya@facebook.com