# OSDF
## A cloud enabled system to store, access, and analyze scientific data

Anup Mahurkar

Institute for Genome Sciences

BOSC 2013

# Motivation

- Work based on NIH Human Microbiome Project Data Analysis and Coordination Center

- Manage large collections of data sets

- Decorate data sets with rich metadata

- Develop a framework that could be reused for other collaborative or multi-center projects

- Ease of development using a language agnostic API

- Scalable and cloud-enabled

# What is OSDF?

- Generic extensible framework for associating data with metadata

- Examples of data
  - Reference Databases, Sequenced Reads, Assemblies, Alignments, Annotation

- Examples of metadata
  - Sequencing platform, Library Preparation, Sequencing Strategy, Assembly method, Alignment method, Alignment reference, Subject information

# What is OSDF?

- Includes a mechanism for
    - Defining data model for elements
    - Reliance on Ontologies, and controlled vocabularies
    - Defining relationships between elements
    - Generic RESTful API for accessing and placing data
    - Domain specific API with Perl/Python/Java bindings
    - Versioning and history
    - Access control

# Technologies Used

- JavaScript Object Notation (JSON) objects for modeling data elements
  - Lightweight data interchange format
  - Allows sparse data in elements
  - Easy to generate and parse
  - Compact and human readable
  - JSON Schema for validation
- CouchDB for storing JSON objects
  - Document-oriented database
  - Apache project
  - RESTful JSON API out of the box
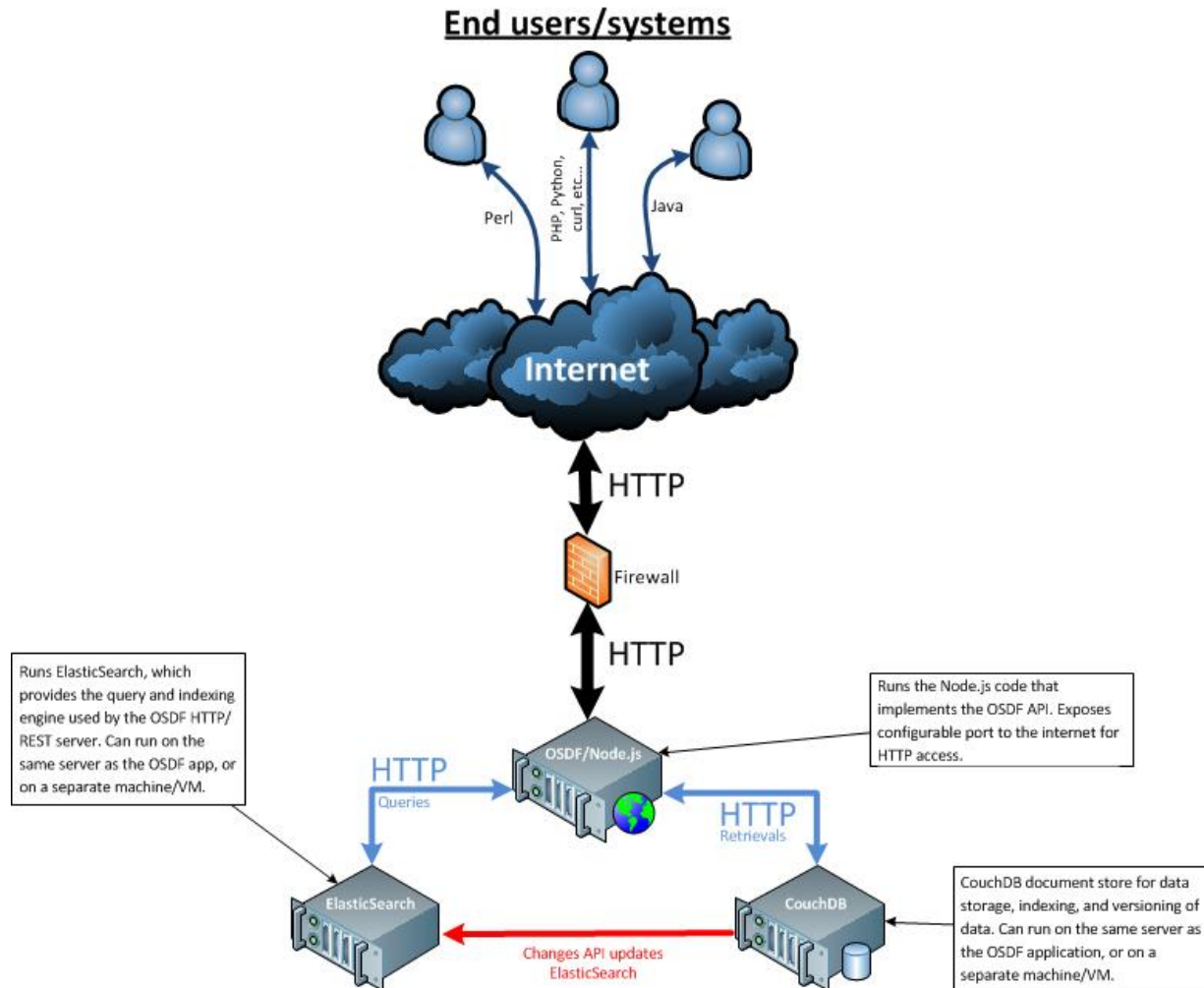  - Real-time bi-directional replication

# Technologies Used

- ElasticSearch
  - Rapid indexing on all keys and attributes of JSON
  - Query language for complex ad-hoc queries
  - Allows wild card, proximity, range, and Boolean operators
- Metadata modeled by CVs, ontologies, standards, and dictionaries
  - MIGS, MIMS, MIMARKS, GO, Relationship
- API implementation using node.js
  - Scalable server optimized for concurrency
  - Implement JSON validation using JSON Schema
- UI Implementation
  - ExtJS, GraphViz, jQuery, D3

# OSDF Architecture

# Schema and Node Example

```
{
  "type": "object",
  "description": "Sample nodes describe physical samples.",
  "properties": {
    "hmp_body_site": {
      "$ref": "body_sites"
    },
    "hmp_supersite": {
      "$ref": "supersites"
    },
    "mimarks_frag": {
      "$ref": "mimarks_frag"
    },
    "mims_frag": {
      "$ref": "mims_frag"
    },
    "visit_number": {
      "type": "integer",
      "minimum": 1,
      "required": true
    },
    "fma_body_site": {
      "title": "Typically a term from the FMA ontology.",
      "type": "string",
      "required": true
    },
    "body_product": {
      "type": "string",
      "required": true
    }
  },
  "additionalProperties": false
}
```

```
{
  "linkage": {
    "part_of": ["c8550ef8d3ea8c9b980650de1c6e86cf"],
    "collected_from": ["c8550ef8d3ea8c9b980650de1c87e0e1"]
  },
  "node_type": "sample",
  "meta": {
    "body_product": "",
    "hmp_supersite": "skin",
    "visit_number": 1,
    "mimarks_frag": {
      "geo_loc_name": "United States of America",
      "samp_mat_process": "",
      "lat_lon": "Unknown",
      "samp_collect_device": "Catch-All sample collection swab",
      "biome": "terrestrial biome [ENVO:00000446]",
      "samp_size": "",
      "feature": "human-associated habitat [ENVO:00009003]",
      "collection_date": "Unknown",
      "env_package": "human-associated",
      "investigation_type": "mimarks-survey",
      "material": "biological product [ENVO:02000043]",
      "rel_to_oxygen": ""
    },
    "fma_body_site": "",
    "hmp_body_site": "left_retroauricular_crease"
  },
  "id": "091f03831014afe9b2da67d698000c51",
  "ver": 1
},
```

# Node Editor

# Current HMP DACC Site

# Programmatic Access

- Generic RESTful API
  - Get
  - Put
  - Update
  - Delete
- Domain Specific API
  - Metagenomics Perl API
    - Sequence data R operations
    - Assembly CR operations
    - Compress files
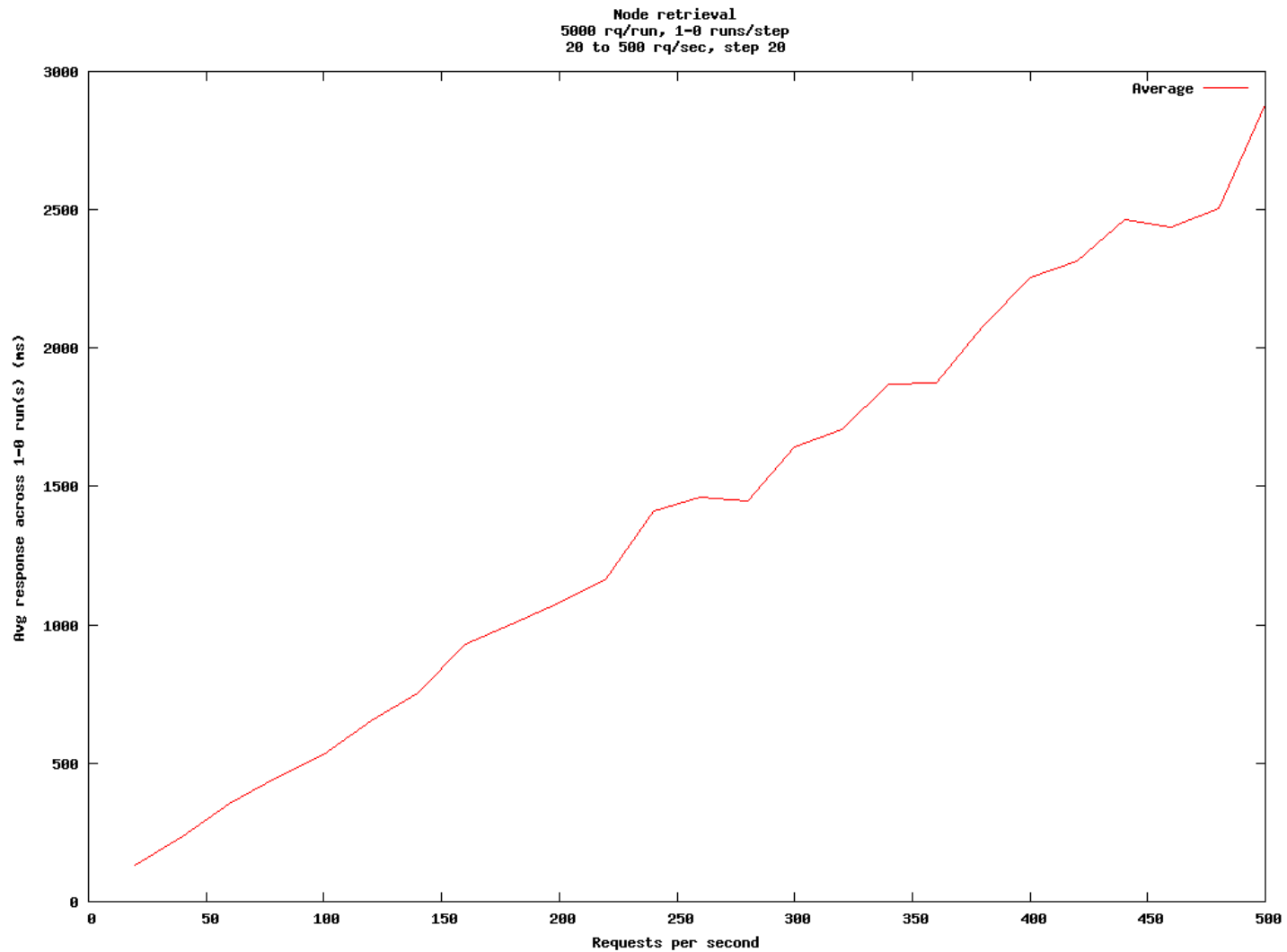    - Upload files
    - Link nodes

# Performance Testing

- Server specifications
  - VM running on VMWare ESX server
  - Ubuntu 12.04
  - 4 vCPUs, 16GB RAM
  - Node.js v0.10
  - Elastic Search v0.90
  - CouchDB v1.2
- Client specifications
  - 8 cores, 6GB RAM
- Benchmark
  - Use 'ab' Apache Server Benchmarking Tool
  - Run 5000 operations with increasing concurrency from 20 to 500 operations
  - { "query": {"term": {"node_type": "sample" }}}

# Performance Numbers – Node Retrieval



Node retrieval
5000 rq/run, 1-0 runs/step
20 to 500 rq/sec, step 20

# Future Development

- Deployment in the cloud

- Replication

- Distributed deployment

- Performance enhancements

- Richer microbiome domain API

- Language bindings in Perl/Python/Java

- Host 1000 Genomes Data

# *Acknowledgements*

- Principal Investigator
  - Owen White

- Institute for Genome Sciences
  - Sonia Agrawal
  - Cesar Arze
  - Jonathan Crabtree
  - Noam Davidovics
  - Victor Felix
  - Aaron Miller
  - Mike Schor

- University of Colorado Boulder
  - Rob Knight

- University of Chicago/Argonne National Lab
  - Folker Meyer
  - Narayan Desai
  - Jared Wilkening

# Resources

- OSDF Website
  http://osdf.igs.umaryland.edu
- OSDF Code
  http://sourceforge.net/projects/osdf/
- Metagenome API site
  http://sourceforge.net/projects/metagenosdf/
- DACC Website
  http://hmpdacc.org

  We invite you to take a look and participate

# Performance Numbers – Node Creation

# Performance Numbers – Node Query



Simple query (all samples)
1000 rq/run, 1 runs/step
20 to 500 rq/sec, step 20