

# HOW WE MADE DATA PROCESSING SCALABLE AT NUGG.AD

- young and energetic
- older and a little wiser

nugg.ad provides targeting service for online advertising

- socio-demographics ( 10 in total )
- product interests ( 48 in total )
- advertisement clicks
- Frequency
- geographic

# online service creates a lot of data

- page clicks
- advertisement clicks
- predictions
- surveys

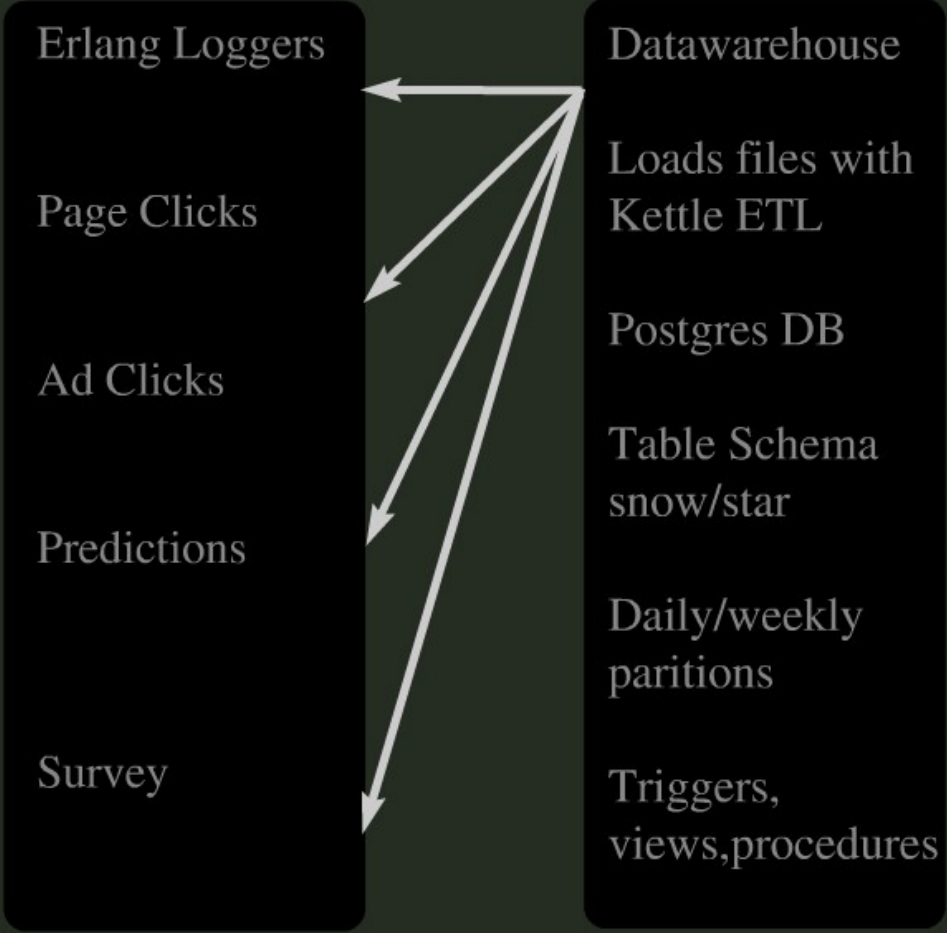
Find a way to use all of this data

# what we need to do

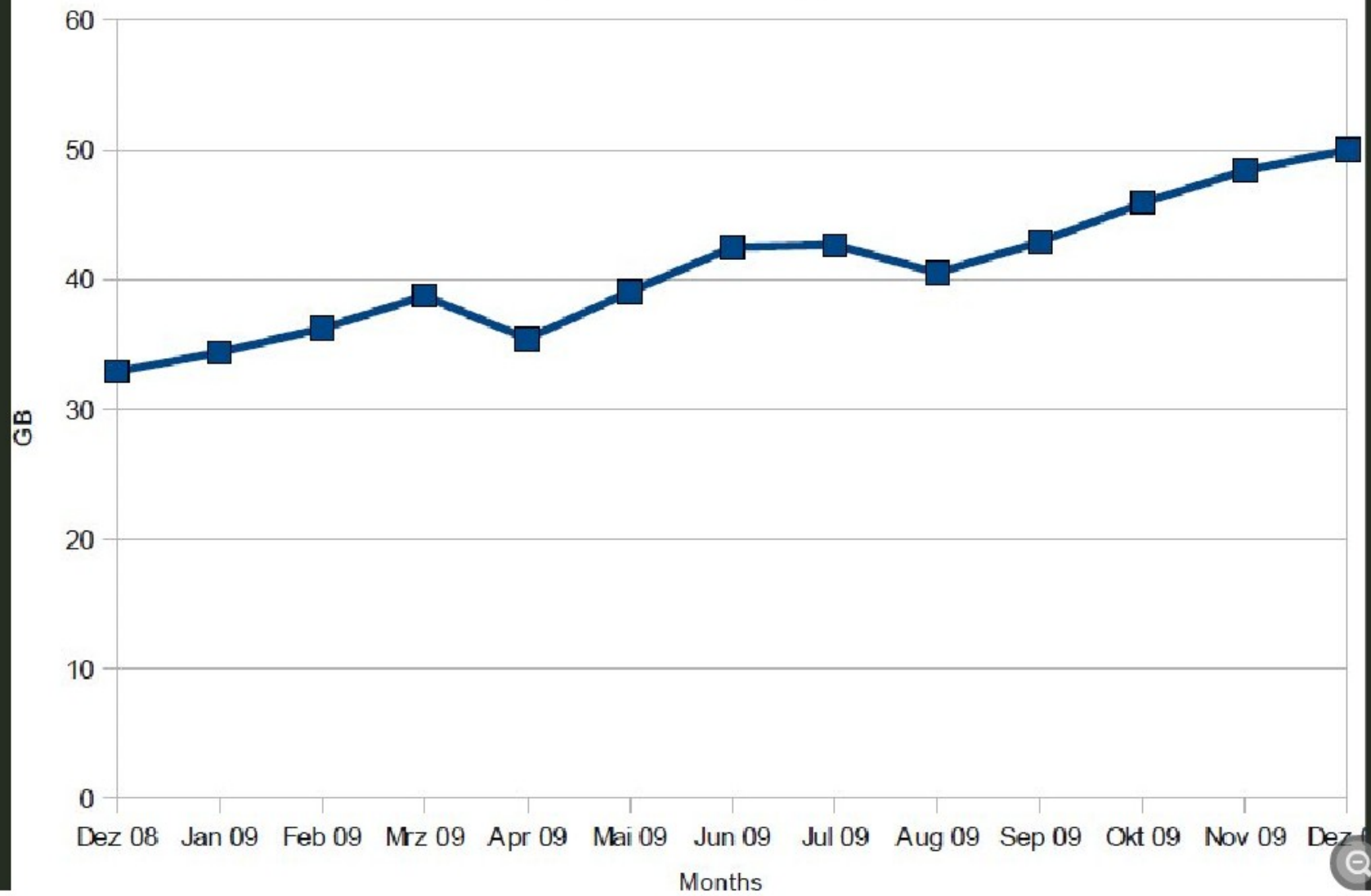
- Aggregate event based logs
- RUN analytics
- create training data
- reports
- create files for db bulk loading

young and energetic

UDP Messages →



### Daily Amount of Log Data





## headache now a migraine

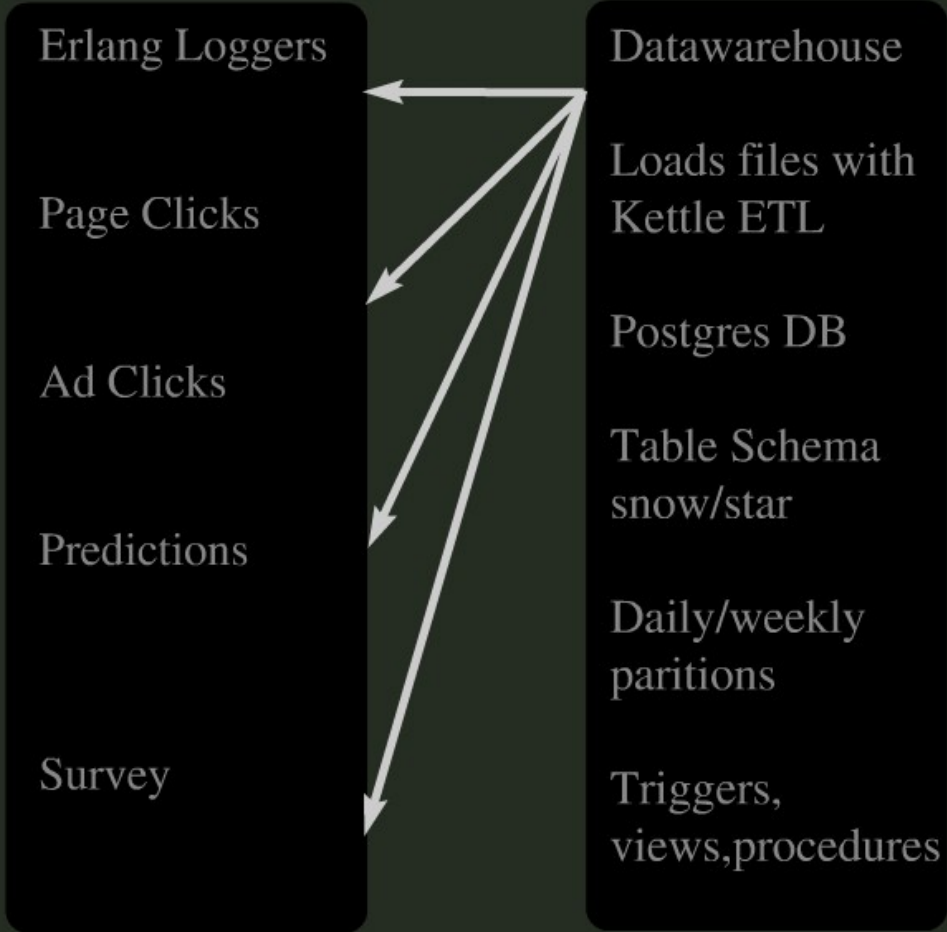
January 2008

12 hours daily data processing  
2 days to create weekly reports  
2 days to create training data  
2 days to create db export files

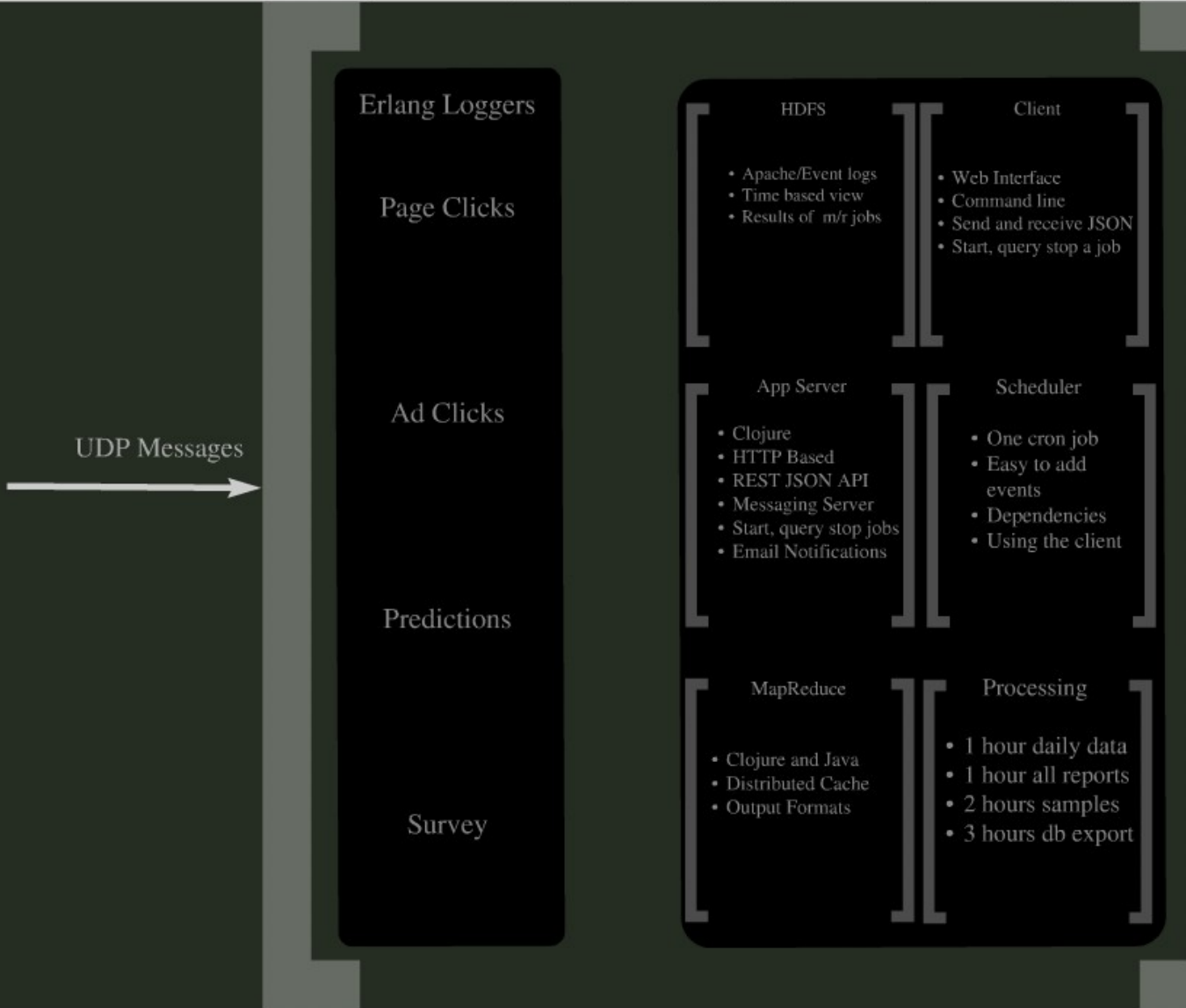
January 2009

- 23 hours daily data processing
- 5 days to create weekly reports
- 4 days to create training data
- 5 days to create db export files

UDP Messages →



older and a little wiser



## App Server

- Clojure
- HTTP Based
- REST JSON API
- Messaging Server
- Start, query stop jobs
- Email Notifications

So

- On
- Eas
- eve
- De
- Us

# HDFS

- Apache/Event logs
- Time based view
- Results of m/r jobs

# Client

- Web Interface
- Command line
- Send and receive JSON
- Start, query stop a job

```
{ "request" :  
  {  
    "action" : "days-click-report",  
    "input" : "/aggregate/click/month",  
    "output" : "/report/click/month"  
  }  
}
```





Erlang Loggers

Page Clicks

Ad Clicks

Predictions

Survey

HDFS

- Apache/Event logs
- Time based view
- Results of m/r jobs

Client

- Web Interface
- Command line
- Send and receive JSON
- Start, query stop a job

App Server

- Clojure
- HTTP Based
- REST JSON API
- Messaging Server
- Start, query stop jobs
- Email Notifications

Scheduler

- One cron job
- Easy to add events
- Dependencies
- Using the client

MapReduce

- Clojure and Java
- Distributed Cache
- Output Formats

Processing

- 1 hour daily data
- 1 hour all reports
- 2 hours samples
- 3 hours db export

Messages



# Scheduler

- One cron job
- Easy to add events
- Dependencies
- Using the client

```
{"events":  
  [  
    {  
      "name"      : "click-report",  
      "position"  : 1,  
      "schedules" : [{"name": "days", "input": "/aggregate/click/month", "output": "/report/click/month"}]  
    }  
  ]  
}
```



```
{ "request" :  
  {  
    "action" : "days-click-report",  
    "input" : "/aggregate/click/month",  
    "output" : "/report/click/month"  
  }  
}
```



Erlang Loggers

Page Clicks

Ad Clicks

Predictions

Survey

HDFS

- Apache/Event logs
- Time based view
- Results of m/r jobs

Client

- Web Interface
- Command line
- Send and receive JSON
- Start, query stop a job

App Server

- Clojure
- HTTP Based
- REST JSON API
- Messaging Server
- Start, query stop jobs
- Email Notifications

Scheduler

- One cron job
- Easy to add events
- Dependencies
- Using the client

MapReduce

- Clojure and Java
- Distributed Cache
- Output Formats

Processing

- 1 hour daily data
- 1 hour all reports
- 2 hours samples
- 3 hours db export

Messages



# MapReduce

- Clojure and Java
- Distributed Cache
- Output Formats

- 1
- 1
- 2
- 3

```
"Map input records possible definitions:
net user page-impressions
net user demographics-prediction primary-key
net user interests-prediction primary-key"

"Reduce output record:
demographic-prediction-primary-key interest-prediction-primary-key page-impression user net"

(defstruct schema :demographics-prediction :interests-prediction :col :user :net)

(def *key-len* 2)

(defn prediction-record?
  [x]
  (contains? x "prediction"))

(defn assoc-prediction
  [schema row]
  (let [name (keyword (.first row))
        value (.last row)]
    (assoc map name value)))

(defn assoc-click
  [schema row]
  (let [pi (.first row)]
    (assoc map :pi pi)))

(defn db-row
  [schema col]
  (reduce (fn [schema val]
            (if (prediction-record? (.first val))
                (assoc-prediction schema val)
                (assoc-click schema val)))
          schema col))

(defn click-prediction-mapper-map
  [this k w out reporter]
  (let [[user net :as parts] (split-line v)]
    (collect out (longPair net user) (sec-tuple (drop *key-len* parts)))))

(defn click-prediction-reducer-reduce
  [this k v out reporter]
  (let [user (.first k)
        net (.last k)]
    (row (vals (db-row (struct-map forecast-schema :user user :net net) (iterator-seq v)))
        (when (no-any? nil? row)
          :collect out (NIL-Writable/get) (sec-tuple row))))))
```

- Users
- Clients
- Page Im
- You sele
- Desktop
- Local file
- User
- File
- App
- Recent
- FileVault
- Preferences
- System
- Keychain
- Time Machine



"Map Input records possible definitions:

```
net user page-impressions
net user demographics-prediction primary-key
net user interests-prediction primary-key"
```

"Reduce output record:

```
demographic-prediction-primary-key interest-prediction-primarty-key page-impression user net"
```

```
(defstruct schema :demographics-prediction :interests-prediction :pi :user :net)
```

```
(def *key-len* 2)
```

```
(defn prediction-record?
```

```
  [x]
  (contains? x "prediction"))
```

```
(defn assoc-prediction
```

```
  [schema row]
```

```
  (let [keys (keyword (first row))
```







Media Selection

Sociodemographics

Product Interests & Channels

Clients 19,000  
Page Impr. 1,096,500

Your selected target group

Sociodemographics

Combination AND

Gender

Male

Age

30 to 39 years

Additional Sociodemographics

Combination AND

employment status

Not working or temporarily unemployed

save

Sociodemographics

Gender

Male  Female

Age

14 to 19 years  20 to 29 years  30 to 39 years  
 40 to 49 years  50 to 59 years  60 years and older

Additional Sociodemographics

Employment status

Fully or partly employed  Pensioner  
 Not working or temporarily unemployed

Level of education

none  primary or lower secondary level (low)  upper secondary level (medium)  
 post-secondary / third level (high)

Personal income

no own income  low income  average income  
 above average  high income  very high income

Responsibility for household

Is responsible for household  Is not responsible for household

Main income-earner in household

Is main income-earner in household  Is not main income-earner in household

Erlang Loggers

Page Clicks

Ad Clicks

Predictions

Survey

HDFS

- Apache/Event logs
- Time based view
- Results of m/r jobs

Client

- Web Interface
- Command line
- Send and receive JSON
- Start, query stop a job

App Server

- Clojure
- HTTP Based
- REST JSON API
- Messaging Server
- Start, query stop jobs
- Email Notifications

Scheduler

- One cron job
- Easy to add events
- Dependencies
- Using the client

MapReduce

- Clojure and Java
- Distributed Cache
- Output Formats

Processing

- 1 hour daily data
- 1 hour all reports
- 2 hours samples
- 3 hours db export

Messages



# Processing

- 1 hour daily data
- 1 hour all reports
- 2 hours samples
- 3 hours db export

## HOW WE MADE DATA PROCESSING SCALABLE AT NUGG.AD

- young and energetic
- older and a little wiser