# Optimizing Your Search Experience

How-To Webinar

# Agenda

- Basic Search Structure
- Setting Search Performance Expectations
- Search Optimization Tools
  - Field Extraction Rules
  - Partitions
  - Scheduled Views
- Demo
- Q&A

# Basic Search Structure

# Search Structure

Keywords and operators (separated by pipes) that build on top of each other

Syntax:

**metadata tags + keywords | parse | filter | aggregate | sort | limit**

Example Search:



```
Unnamed Search        +

_sourceCategory=Apache/Access and GET
| parse "GET * HTTP/1.1\" * * \"*\"" as url, status_code, size, referrer
| where !(status_code = 200 and status_code=304)
| count by status_code
| sort by status_code asc
| limit 10
```

**Results**

# Metadata Fields

- All messages are tagged during data ingest
- Metadata fields are configured as part of Collector and Source setup

| Name | Description |
|------|-------------|
| _collector | Name of collector when installed |
| _source | Name of the source defined during configuration |
| _sourceHost | The host name of the source |
| _sourceName | The name of the log file (including path) |
| _sourceCategory | Category associated with the source |

- Properly categorizing your data leads to more efficient searches
- Good Source Category, Bad Source Category

# Keyword Search

- Case Insensitive unless string is in double quotes

- Wildcard Support (e.g. ERR*)

- Boolean Logic Support
  - AND
  - OR
  - !(A OR B)

- Combine keywords with metadata fields for the best performance

- Bloom filters
  - Using keywords helps bloom filters retrieve data very quickly

# Processing Your Search Request

**Initiate**
- Queries are rewritten automagically
- The Sumo Logic service calls backend clusters to kickoff the request

**Reduce**
- Sumo locates indices that contain data for search time-range
- Bloom filters further eliminate indices where keywords are not contained

**Data Retrieval**
- Everything through the first pipe is retrieved
- Data is carried forward

**Parallelize**
- Remaining operations are conducted
- If aggregation is involved, we look for opportunities to parallelize the operation

# Develop Good Search Habits

- Use metadata and keyword combinations to reduce scope

- Shorten your time-range down as much as possible

- Limit result sets before aggregating data → where user=a | count by user

- Use parse anchor instead of parse regex for structured messages

- Avoid the use of expensive parse regex tokens like .* → \d{2,10}

- Add line breaks after each operation

- Use pre-extracted fields where possible (to be discussed later)

# Search Performance Expectations

# The Time Range Effect

- More recent data can be accessed quickly
    - We do something special when scanning the last 24 hours of events
    - Why? Over 90% of searches are executed against recent data

| Last 24 Hours | 🕐 | Start |

☐ Use Receipt Time

- Test queries on very recent data first before saving or publishing
- Our main performance metric (speedup) is essentially a ratio that divides the time-range used by the time it takes for data to return.

# Review Your Data Source Time Zone Settings

- Leads to a large gap between message time and receipt time

- Causes backend fragmentation and can affect search speed

- Support of Java 6 Time Zone formats
  - Pacific Standard Time; PST; GMT-08:00
  - -0800
  - **NOT** US/Pacific

- Data integrity will be questioned by users

- Knowledge Base Article: Large Time Discrepancies
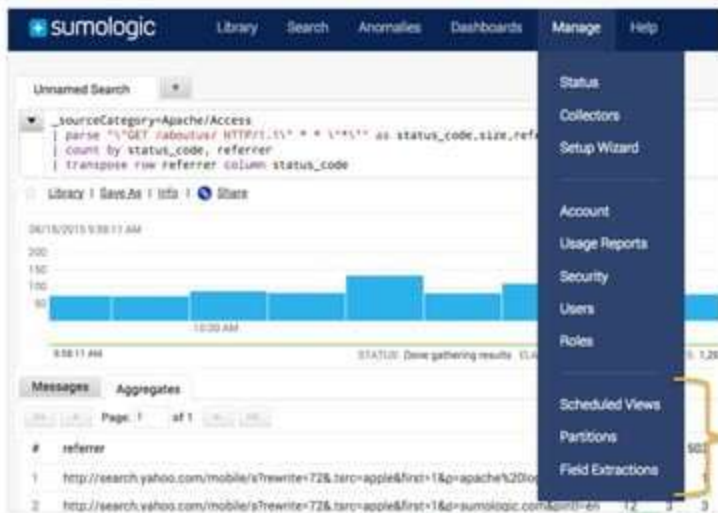
# Compute-Intensive Operations

- Multiple .* tokens in a single parse regex statement
- Parse using public library (apache/access, iis, cisco/asa, windows/2008)
  - Try to borrow from Field Extraction Rule templates
- LogCompare and LogReduce
- Join
  - Time to run exponentially increases when extending your time-range
- Transaction, Transactionize and Merge
  - Try and limit the 'timewindow' parameter for finding corresponding events
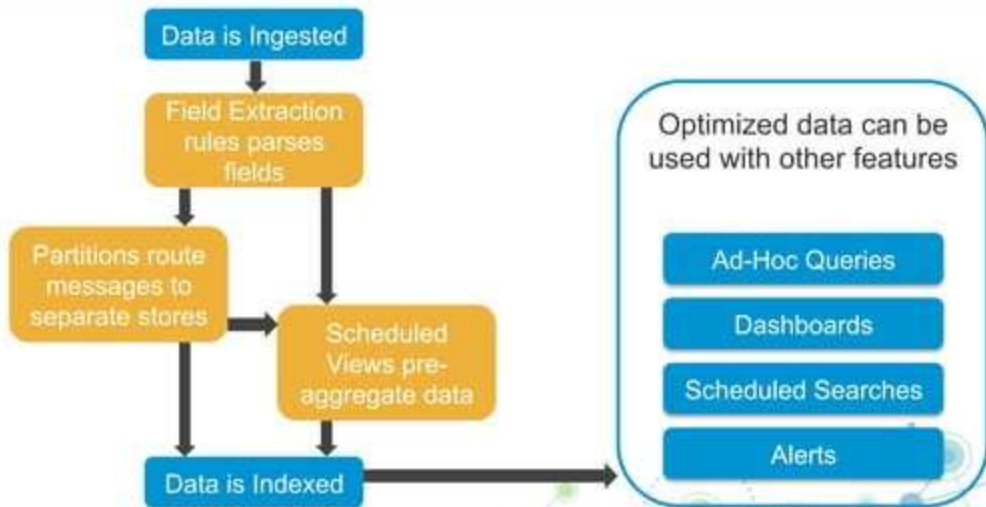- Outlier / Predict

Performance Optimization Tools

# Managing Search Optimization Tools

# How Data is Optimized for Search
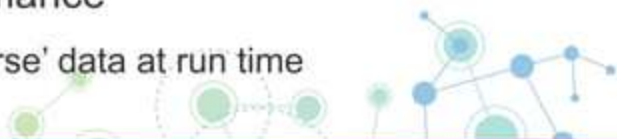
# Benefits of Field Extraction Rules

- Extract fields at the time of ingest
- Standardize Searches and Field Names for users
- Simplify searches
  - Narrow results within search scope instead of using 'where' operator (e.g. _sourceCategory=nginx status_code=404)
- Improves Search Performance
  - Eliminates the need to 'parse' data at run time

# When to Use Field Extraction Rules

- The same (or very similar) parse statement is being used over and over

- Parsing over a large volume of data

- Constantly filtering data based off of parsed fields

- Disparate logs need to be joined using a Unique ID
    - Session ID
    - User Name
    - Process ID

- Syslog Metadata Overrides

# Create Field Extraction Rule



**Edit Field Extraction Rule**

This form allows you to edit a field extraction rule. Enter a name, scope, and fields. You may also choose from a list of parse expression templates instead of creating your own.

**Rule Name** * — Apache Access Log

**Scope** * — _sourceCategory = Apache/Access

*Use Scope to define what data this FER applies to*

**Parse Expression** * — parse regex "*(?<src_ip>\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3})" | parse regex "(?<method>[A-Z]+)\s(?<url>\S+)\sHTTP/(\d\.)+\"\s(?<status_code>\d+)\s(?<size>[\d-]+)\s\"(?<referrer>.*?)\"\s\"(?<user_agent>.*?)\".*"

*Use Regex to create your parse expression*

**Templates** — ✓ Select a Parse Template (optional)
- Akamai Cloud Monitor
- Apache Access Logs
- Apache Tomcat Access Logs
- AWS Cloud Trail Logs
- AWS Elastic Load Balancing Logs
- AWS S3 Usage Logs
- Microsoft IIS Logs
- Nginx Logs
- Palo Alto Networks
- Varnish Logs

*Templates exist for common sources*

Cancel    Save As

# Using Pre-Parsed Fields When Querying



_sourceCategory=Apache/Access method=GET

> With FERs, parsed fields are available to use in your keyword search

> Parsed fields are available in your Field Browser for further analysis

## Field Extraction Rule Recommendations

- Test the rule by running a search over a small time-range that has data

- The scope and parse statement should not change

- Ensure your rule covers common searches

- Only extract the minimum fields necessary

  - Use 'fields' operator to limit results

# FER Caveats

- Max of 50 Rules
- Max of 200 Total Fields
- Supported Operators
  - Parse Anchor / Regex / Nodrop
  - Double
  - Fields
  - Num
  - If
  - Where
  - Concat
  - Keyvalue (not 'kv auto')
  - **NEW!** JSON (not 'json auto')

**NOTE**: Deleted rules and fields defined in them will still count towards the max

# Partitions

## Benefits of Partitions

- Divides your data into smaller chunks to be searched on

- Takes advantage of your source categorization; similar data can actually be grouped together

- Improves performance when used in searches

- It can eliminate the need for lengthy scope definitions

## When to Create Partitions

- Sets of data are being searched in isolation
- A large amount of data being sent daily (> 5 GB's)
  - Navigate to Manage → Account if you don't know
- Different groups are focused on specific logs
- RBAC filtering is required for data provisioning

# Use Data Volume Index

- Helps to determine possible ways to partition data

- Recommended partition size → Up to 30% of data volume



- Manage → Account → Data Management

- Library → Apps → Data Volume

## Partitions Caveats

- Overlapping data between partitions are counted towards your contracted data volume quota

- Maximum of 500 indexes can be created with no available overrides

- Data cannot be backfilled

- Not editable after creation

# Scheduled Views

# Benefits of Scheduled Views

- Similar to relational DB materialized views

- Allows you to pre-aggregate data

- Allows for long term trending analysis

- Can significantly improve performance for high selectivity queries

    - (_source=A or _source=B) and _sourceName=C and keyword1 and keyword2

- Unlike partitions, data can be backfilled

# When to Use Scheduled Views

- Specific aggregate operators are used heavily in queries
  - Count
  - Sum
- Data is being trended over a long period of time (e.g. Last 30 Days)
  - Failed logins on critical servers
  - Number of 404 errors
- A highly selective query does not perform well

# Scheduled Views Recommendations



- ➕ Include aggregation

- ➕ Timeslice 1m

- ➕ Use queries that are not likely to change

- ➕ Take advantage of existing partitions and FER's

- ➕ Only backfill data needed for analyses

## Scheduled View Caveats

- Data in scheduled views are counted towards your quota
- Parsed fields in views count towards field extraction limitation (200)
- Data can only be backfilled through your plan's retention period
- Not editable after creation
- Supported aggregate operators
    - Difference
    - Count
    - Sum

# Quick Review

# Review: Factors in Search Performance

- Query structure
    - Data Selectivity (keywords, metadata fields)
    - Heavy Operations (join, transaction, summarize)
- Search Time Range
- Possible Time Zone Misconfiguration at Source Level
- Total Data Volume for Account
- Use of Performance Optimization Tools
- Service Anomalies

# Review: Search Optimization Tools

| What I want to do is | Partition | Scheduled View | Field Extraction |
|---|---|---|---|
| Parse the same type of log message repeatedly | | | ✔ |
| Identify long-term trends | | ✔ | |
| Group related data together | ✔ | | |
| Pre-compute or aggregate data before querying | | ✔ | |
| Use RBAC to deny or grant access to the data | ✔ | | |

customer-success@sumologic.com

optimize@sumologic.com