

The ultimate guide for Elasticsearch plugins



Itamar Syn-Hershko

<http://code972.com>

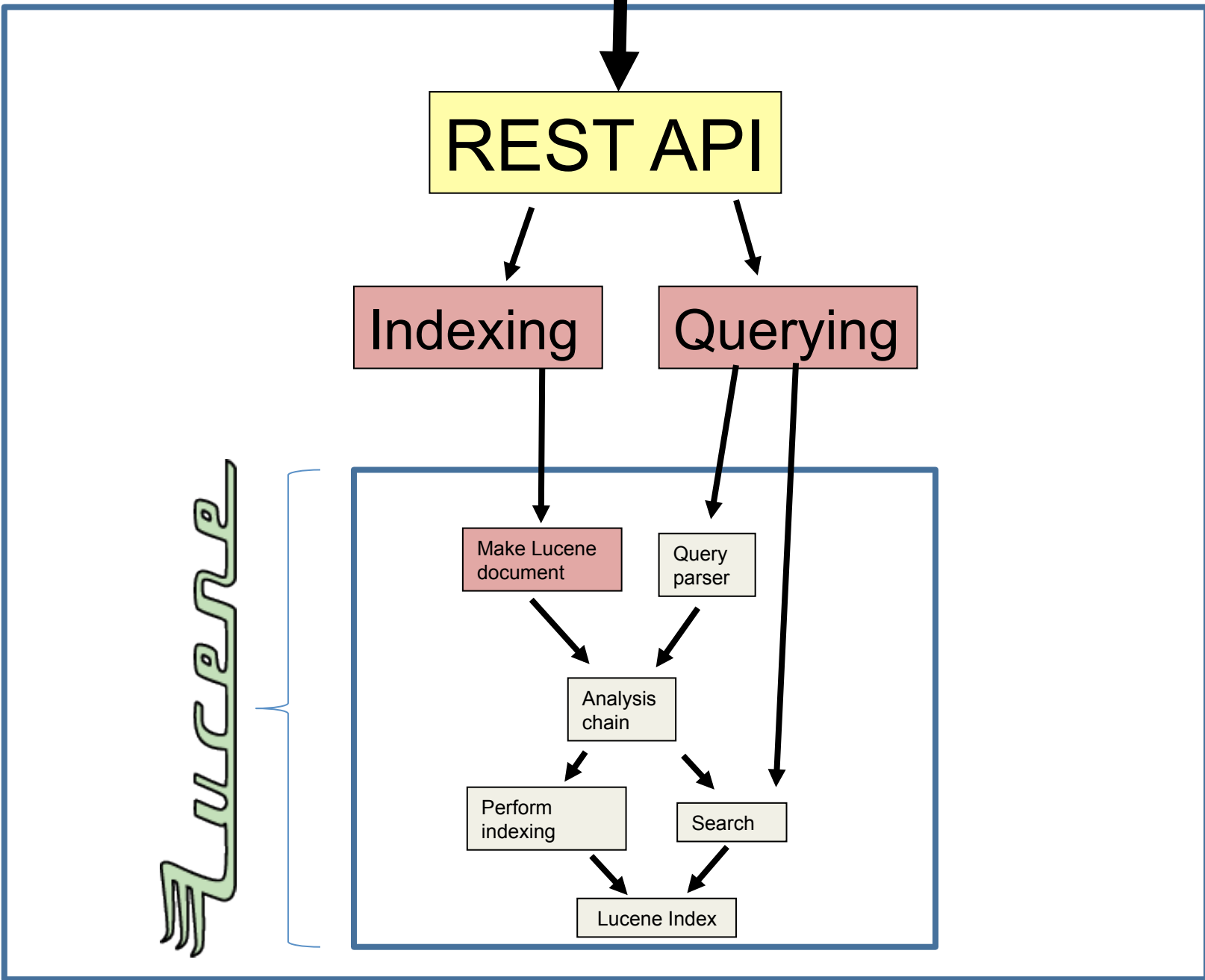
@synhershko

Agenda

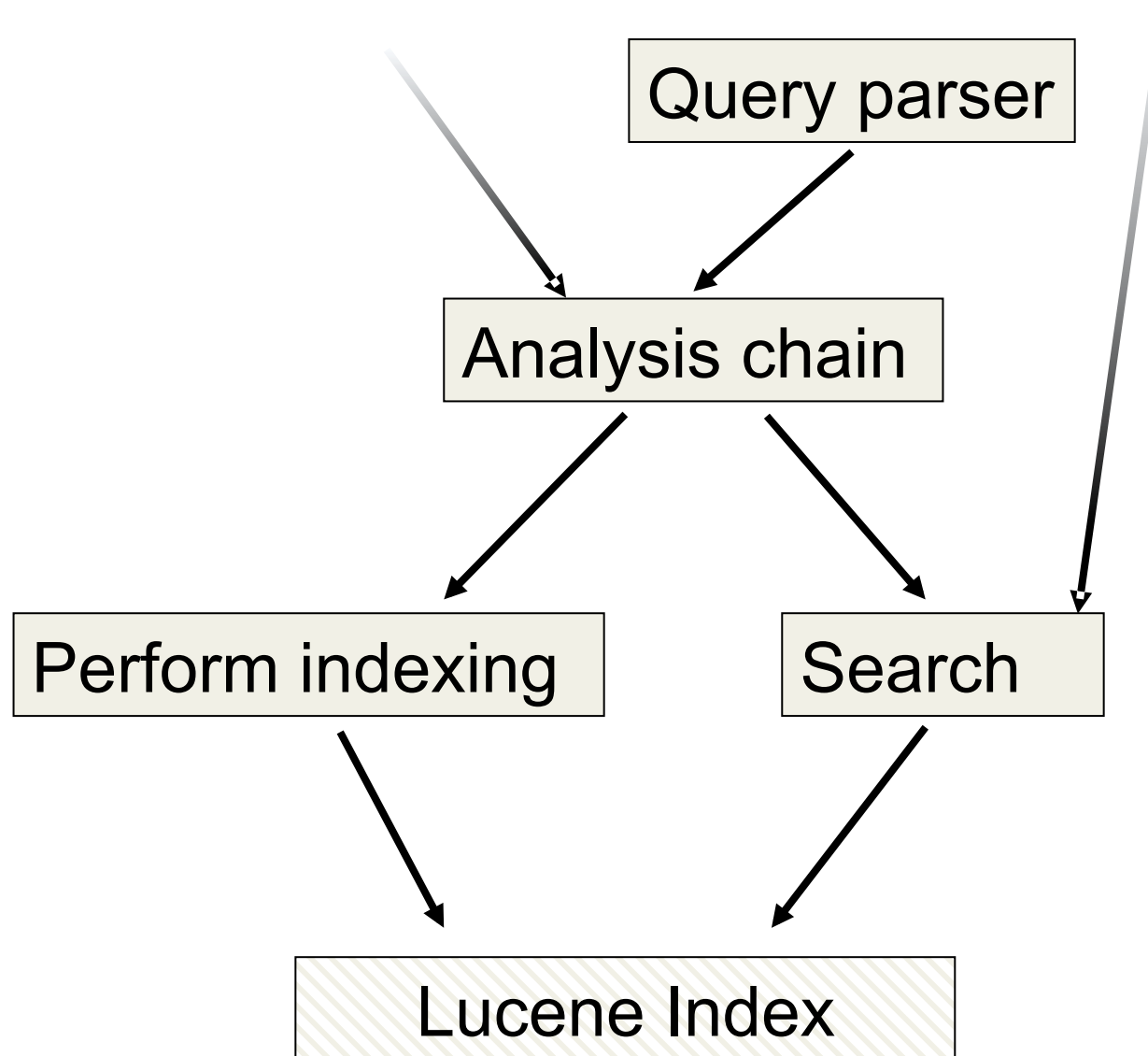
- Integration points & plugin types
- Showcases
- Gotchas
- When-to, How-to
- Q & A

Elasticsearch Server

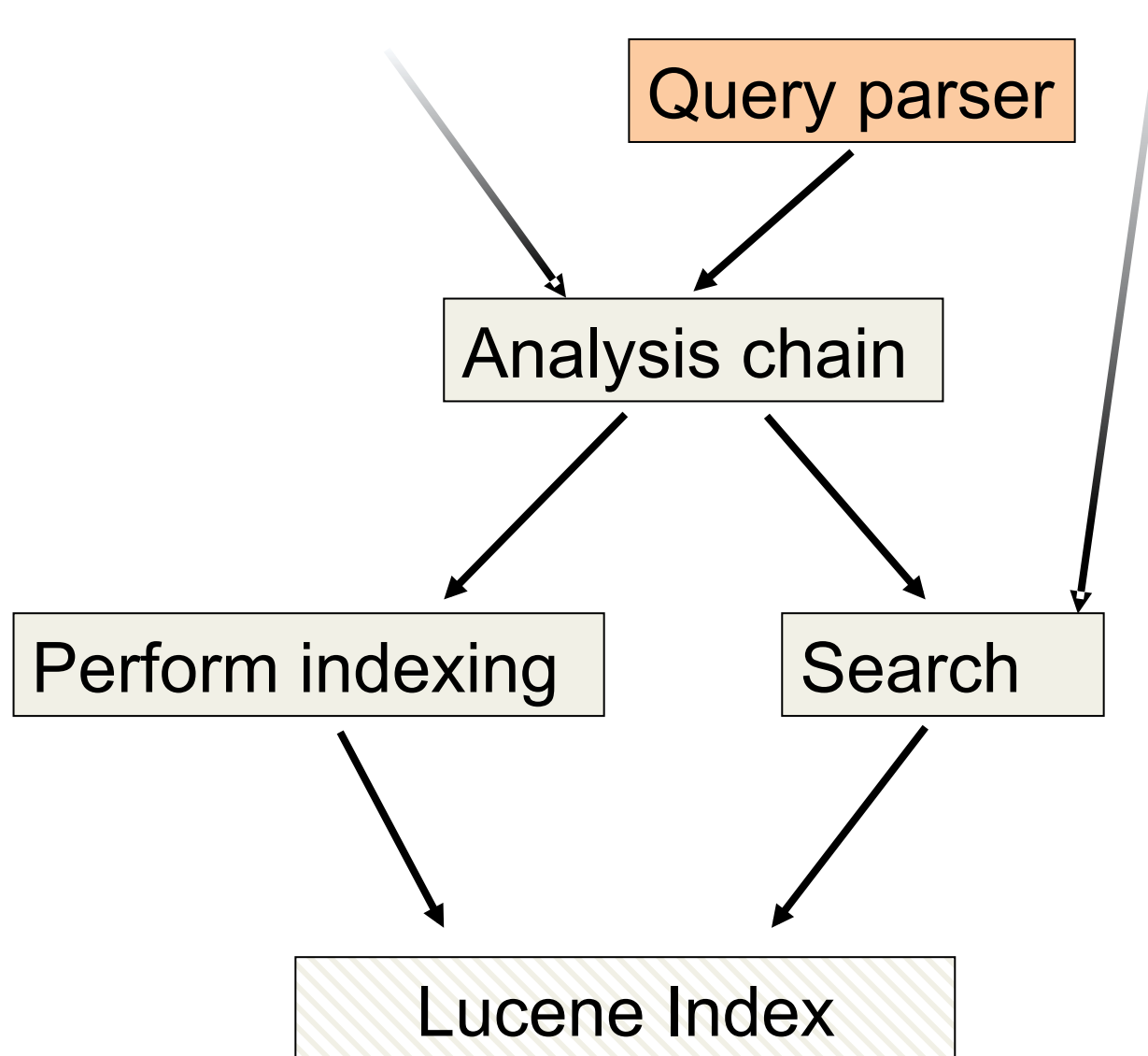
Lucene



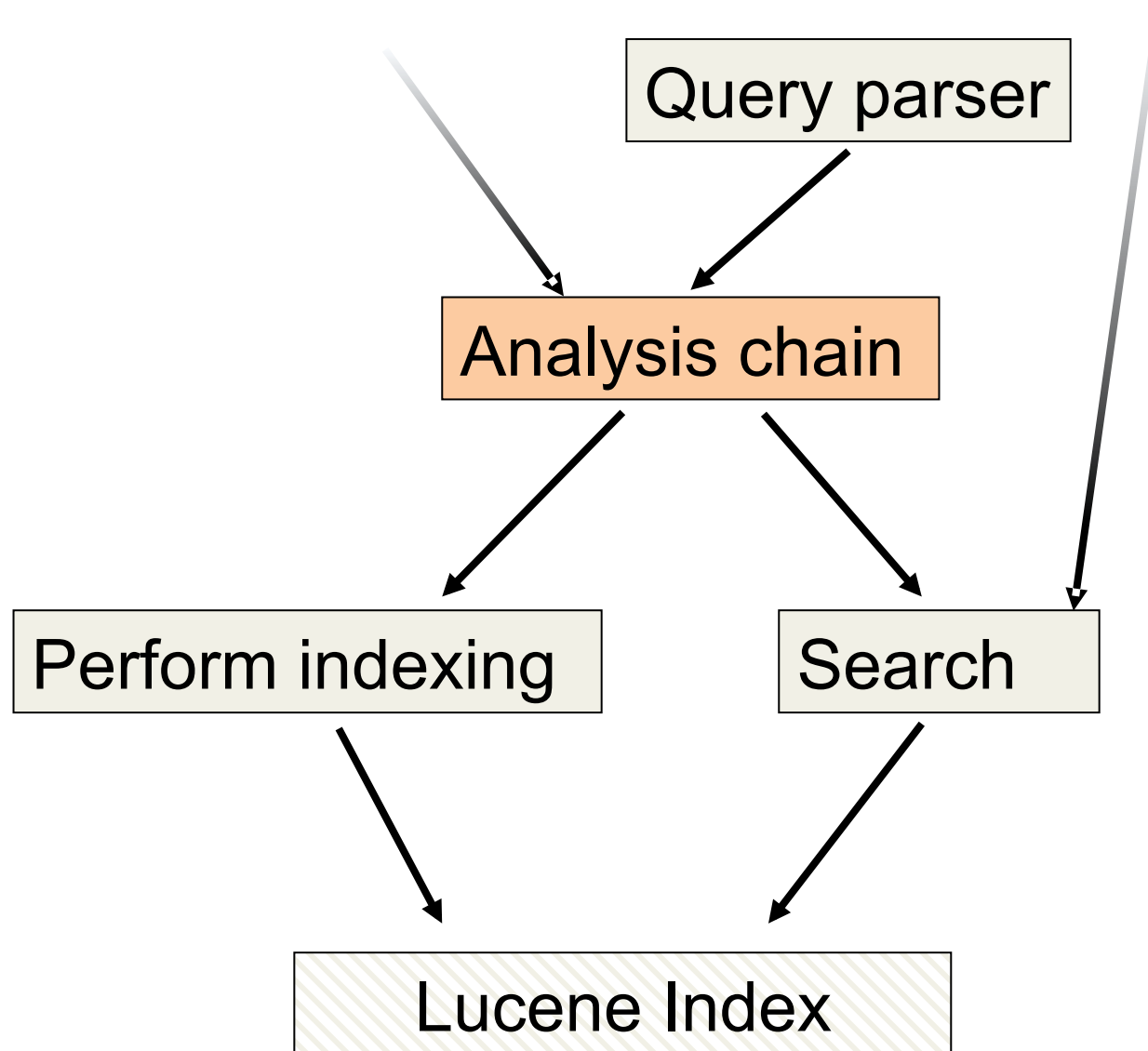
Lucene extension points



Lucene extension points



Lucene extension points



Harry Potter and the Goblet of Fire

Tokenizer

Harry Potter and the Goblet of Fire
Potter the of

Lower case filter

harry potter and the goblet of fire

Stop-words filter

harry potter goblet fire

Step 1: Tokenization

Step 2: Filtering

Welcome to Malmö!

Tokenizer

Welcome Malmö
to

ASCII folding
filter

Welcome Malmo
to

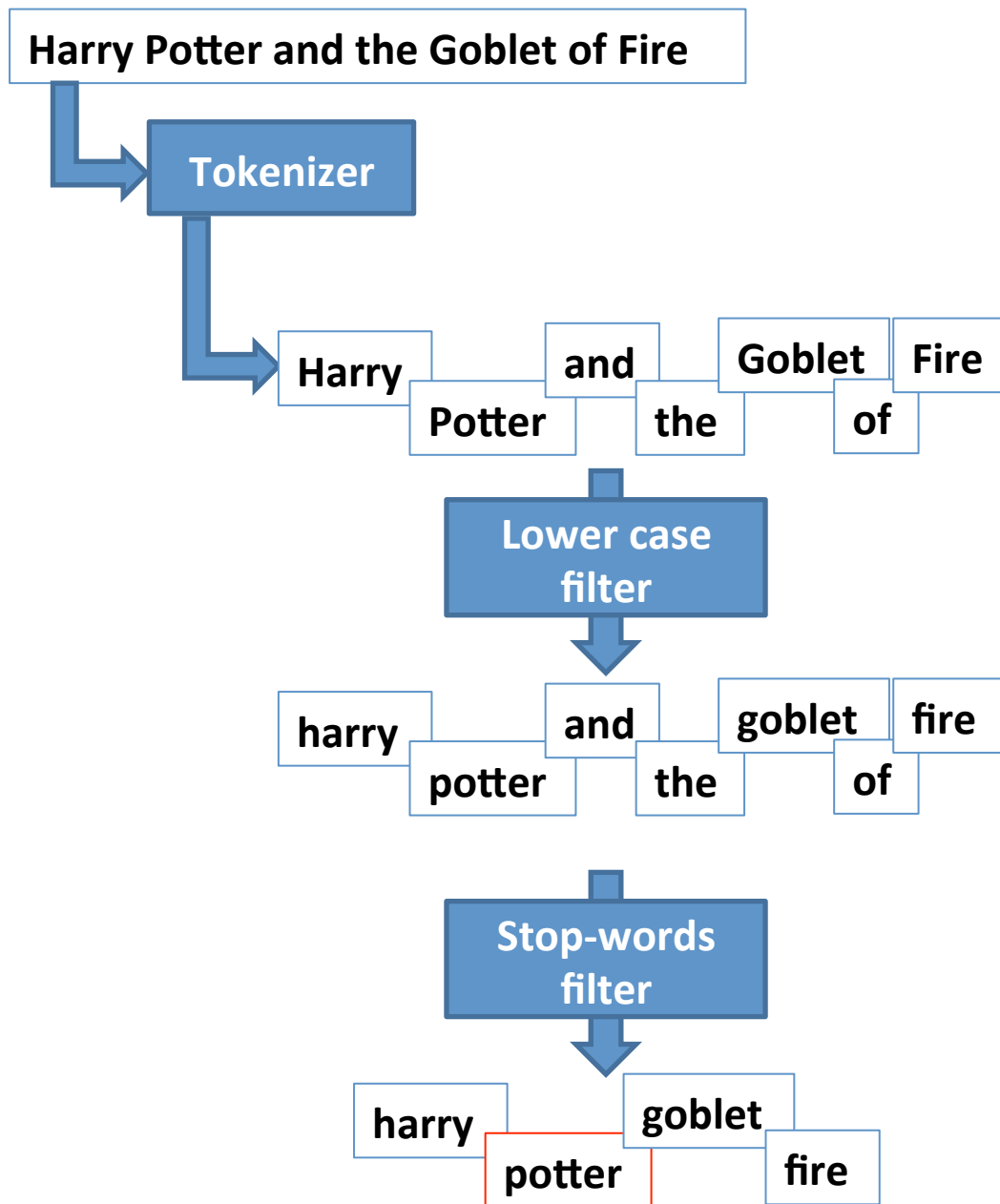
Lowercase
filter

welcome malmo
to

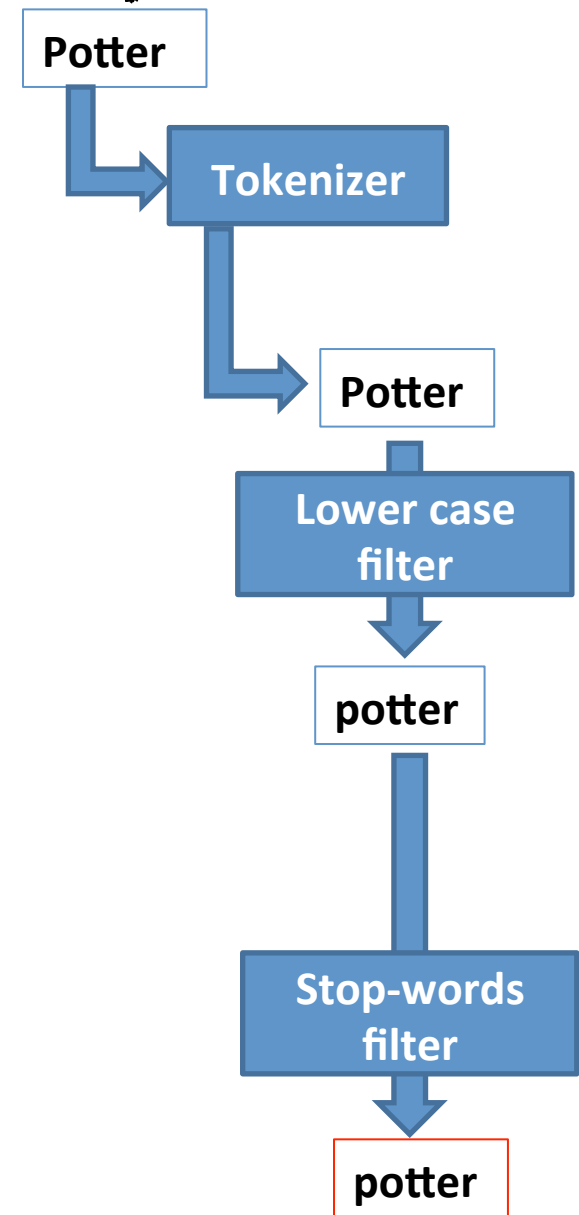
Step 1: Tokenization

Step 2: Filtering

Indexing



Query



itamar@code972.com

Tokenizer

itamar code 972 com

Lower case filter

itamar code 972 com

Step 1: Tokenization

Step 2: Filtering

Try searching on German compound words...

Donaudampfschiffahrtselektrizitätenhauptbetriebswerkbauunterbeamtengesellschaft

Company

The compound word is an example of the virtually unlimited compounding of nouns that is possible in many Germanic languages. [Wikipedia](#)

Feedback

Donaudampfschiffahrtselektrizitäten

Company

The compound word is an example of the virtually unlimited compounding of nouns that is possible in many Germanic languages. [Wikipedia](#)

Analyzers

The quick brown fox jumped over the lazy dog,
bob@hotmail.com 123432.



StandardAnalyzer:

[quick] [brown] [fox] [jumped] [over] [lazy] [dog] [bob@hotmail.com] [123432]

StopAnalyzer:

[quick] [brown] [fox] [jumped] [over] [lazy] [dog] [bob] [hotmail] [com]

SimpleAnalyzer:

[the] [quick] [brown] [fox] [jumped] [over] [the] [lazy] [dog] [bob] [hotmail]
[com]

WhitespaceAnalyzer:

[The] [quick] [brown] [fox] [jumped] [over] [the] [lazy] [dog,]
[bob@hotmail.com] [123432.]

KeywordAnalyzer:

[The quick brown fox jumped over the lazy dog, bob@hotmail.com 123432.]

Custom analyzers from code

```
"index" : {  
  "analysis" : {  
    "analyzer" : {  
      "default" : {  
        "tokenizer" : "standard",  
        "filter" : ["standard", "my_ascii_folding"]  
      }  
    },  
    "filter" : {  
      "my_ascii_folding" : {  
        "type" : "asciifolding",  
        "preserve_original" : true  
      }  
    }  
  }  
}
```


New in Elasticsearch v1.1.0

Showcase: Custom Analyzer - Hebrew analysis plugin for Elasticsearch

- <https://github.com/synhershko/elasticsearch-analysis-hebrew>
- Available on QBox.io

<input type="radio"/>	4	7.5 GB	80 GB ⚡
<input type="radio"/>	8	15 GB	160 GB ⚡
<input type="radio"/>	16	30 GB	320 GB ⚡
<input type="radio"/>	32	60 GB	640 GB ⚡

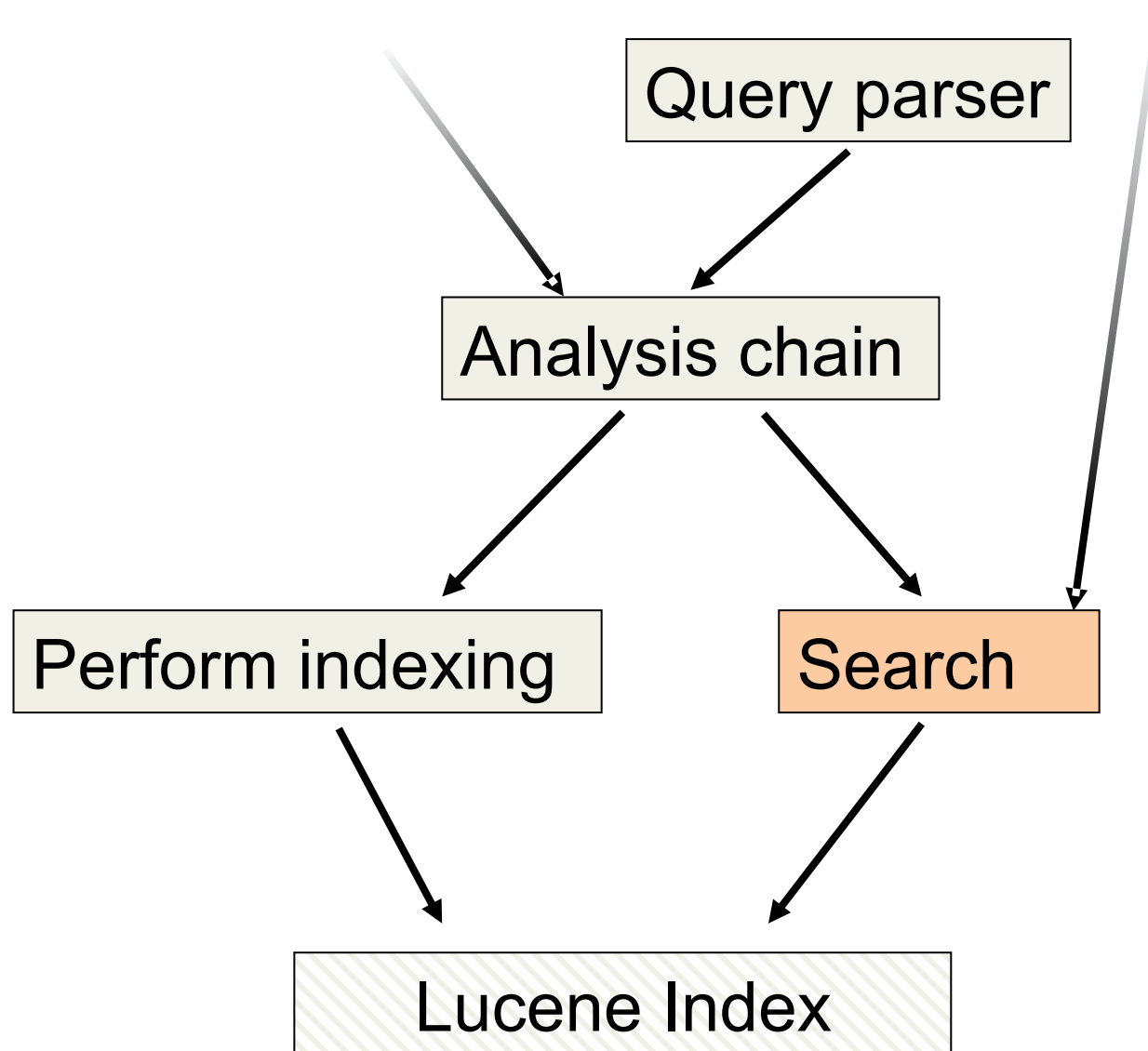
\$0.12/hr (\$84.24/mo)

 1 or 2-node clusters on large boxes can support intensive queries on larger data sets, and several-node clusters on small boxes can support heavy search volumes on smaller data sets.

Options	
<input type="checkbox"/>	Whitelisted IPs
<input type="checkbox"/>	HTTP Basic authentication
<input checked="" type="checkbox"/>	Plugins
<input type="checkbox"/>	Kibana

Plugins	
<input type="checkbox"/>	ICU Analysis
<input type="checkbox"/>	Japanese (Kuromoji) Analysis
<input type="checkbox"/>	Smart Chinese Analysis
<input type="checkbox"/>	Stempel (Polish) Analysis
<input type="checkbox"/>	Combo Analysis
<input checked="" type="checkbox"/>	Hebrew Analysis
<input type="checkbox"/>	CouchDB River

Lucene extension points



Scripting

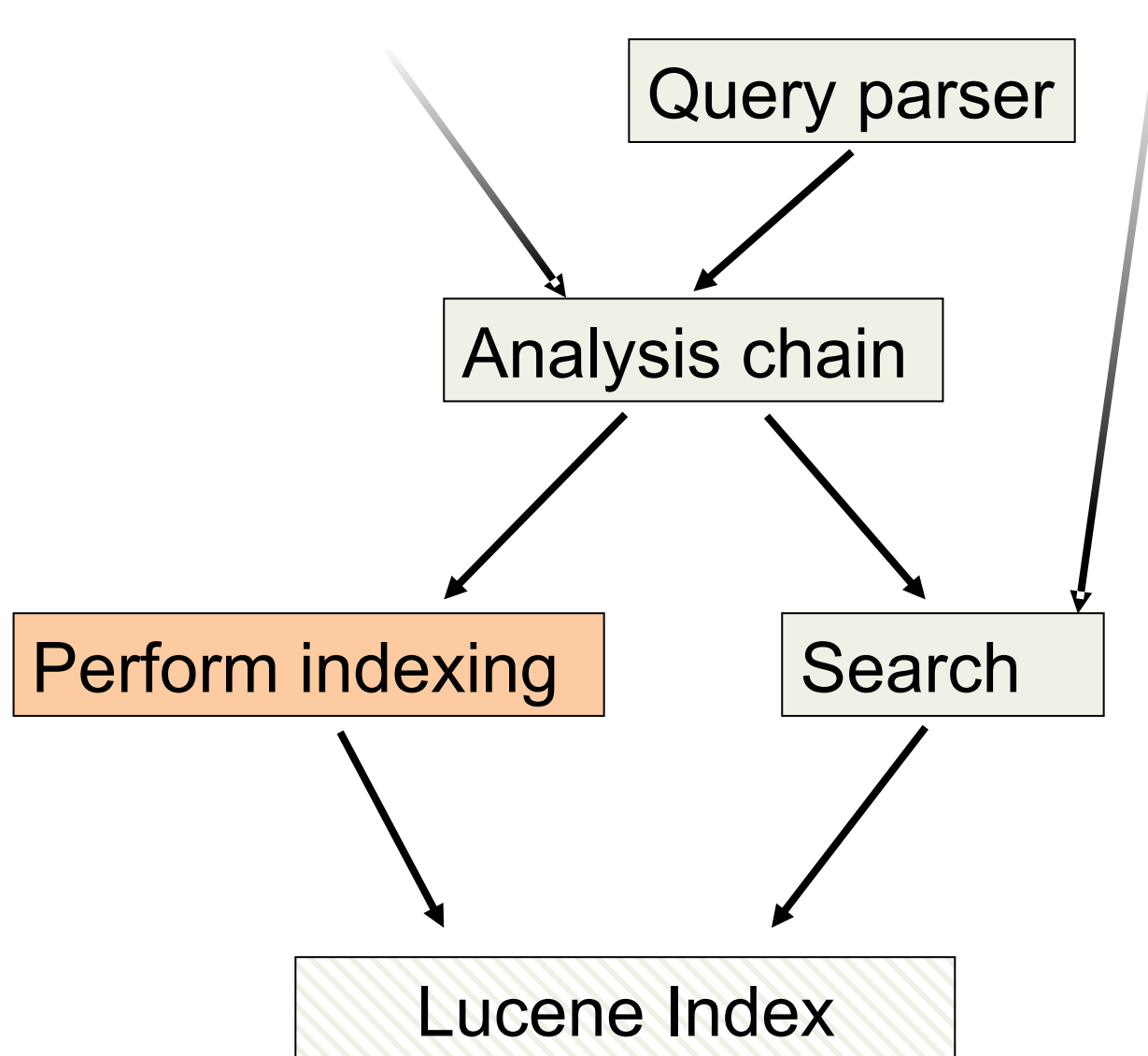
- Sorting, filters, facets, script fields, custom scoring, aggregations, document updates
- MVEL, but others are supported
- Generally speaking: SLOOOOOOOW
- Mostly useful as quick mocks / PoC
- Native scripts using Java by implementing `AbstractExecutableScript` & `AbstractSearchScript`

Custom scoring & similarity

- Function score query
 - Previously known as Custom Score Query
- Similarity

```
"similarity" : {  
  "my_similarity" : {  
    "type" : "DFR",  
    "basic_model" : "g",  
    "after_effect" : "1",  
    "normalization" : "h2",  
    "normalization.h2.c" : "3.0"  
  }  
}
```

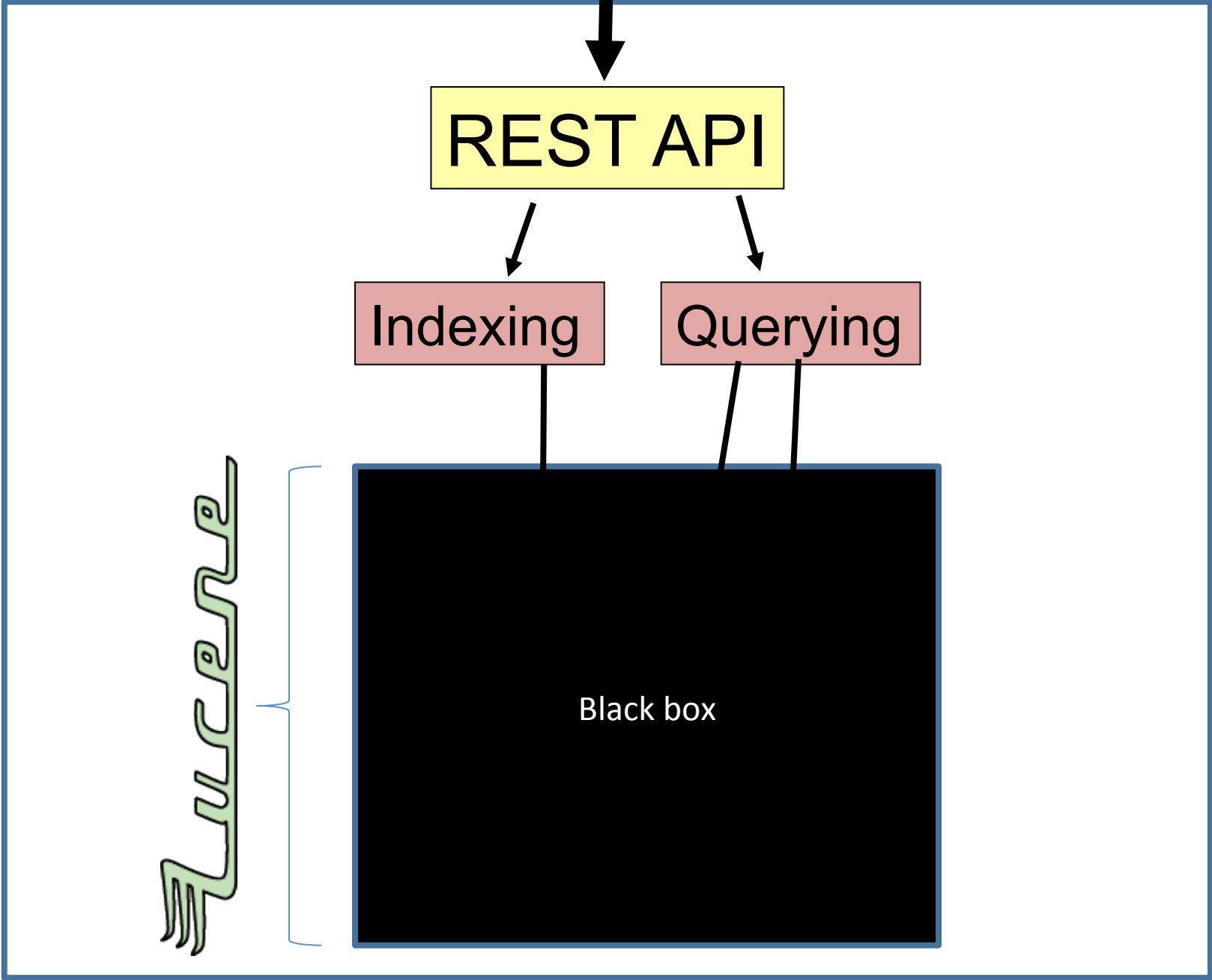
Lucene extension points



Codecs

```
curl -XPUT 'http://localhost:9200/twitter/' -d '{
  "settings" : {
    "index" : {
      "codec" : {
        "postings_format" : {
          "my_format" : {
            "type" : "pulsing",
            "freq_cut_off" : "5"
          }
        }
      }
    }
  }
}'
```

Elasticsearch Server



Controlling shard allocation

- Filtering built in
 - By tags, groups, racks, IPs
 - Black list / white list
- Total shards per node
- Disk based
- EXPERT: Roll your own by implementing AllocationDecider

Custom REST endpoints

```
public class HelloRestHandler implements RestHandler {
    @Inject
    public HelloRestHandler(RestController restController) {
        restController.registerHandler(GET, "/_hello", this);
    }

    @Override
    public void handleRequest(final RestRequest request, final RestChannel channel) {
        String who = request.param("who");
        String whoSafe = (who!=null) ? who : "world";
        channel.sendResponse(new StringRestResponse(OK, "Hello, " + whoSafe + "!"));
    }
}
```

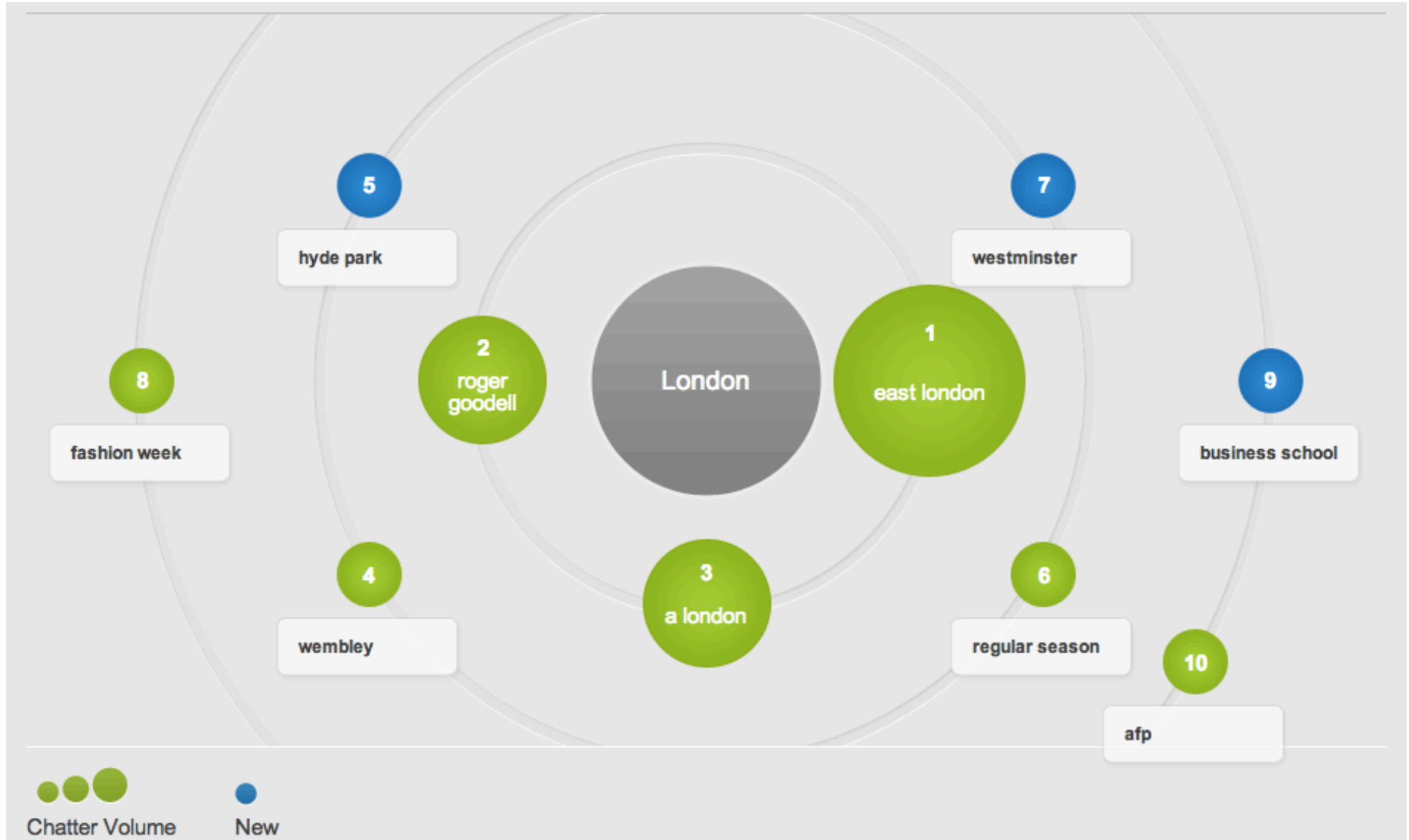
Transports

- Exposes the Elasticsearch RESTful API over protocols other than HTTP
 - Apache Thrift
 - Memcached
 - Servlet
 - Redis
 - ZeroMq

Showcase: Custom percolator



Showcase: The bubble plugin



Site plugins

- Monitoring
 - BigDesk, ElasticHQ, Paramedic, ...
- Hammer (GUI for REST interface)
- Inquisitor (debugging queries)
- SegmentSpy
- WhatsOn

Discovery

- Default is Zen discovery
 - Unicast: I know who my nodes are
 - Multicast: Auto discovery for nodes
- Multicast discovery support for cloud environments
 - AWS
 - Azure
 - Google Compute
- ProTip: Unicast in production unless you know what you're doing
- ZooKeeper plugin

Snapshot / restore repositories

- File system
- AWS S3
- HDFS
- Azure
- Roll your own (e.g. Glacier)

River plugins

- **Obsolete**
- Use the “shoveller” approach
- logstash, stream2es

Summary: Plugin types

- Lucene components
 - Analysis
 - Similarity
 - Scoring
- REST endpoints
- Scripting
- ES infrastructure (Discovery, Transport, Snapshot/restore)
- Site plugins
- ~~River plugins~~

Installing plugins

- Manual under /plugins
- Official / GitHub / Maven installation:

```
plugin --install <org>/<user/component>/<version>
```

- From zip:

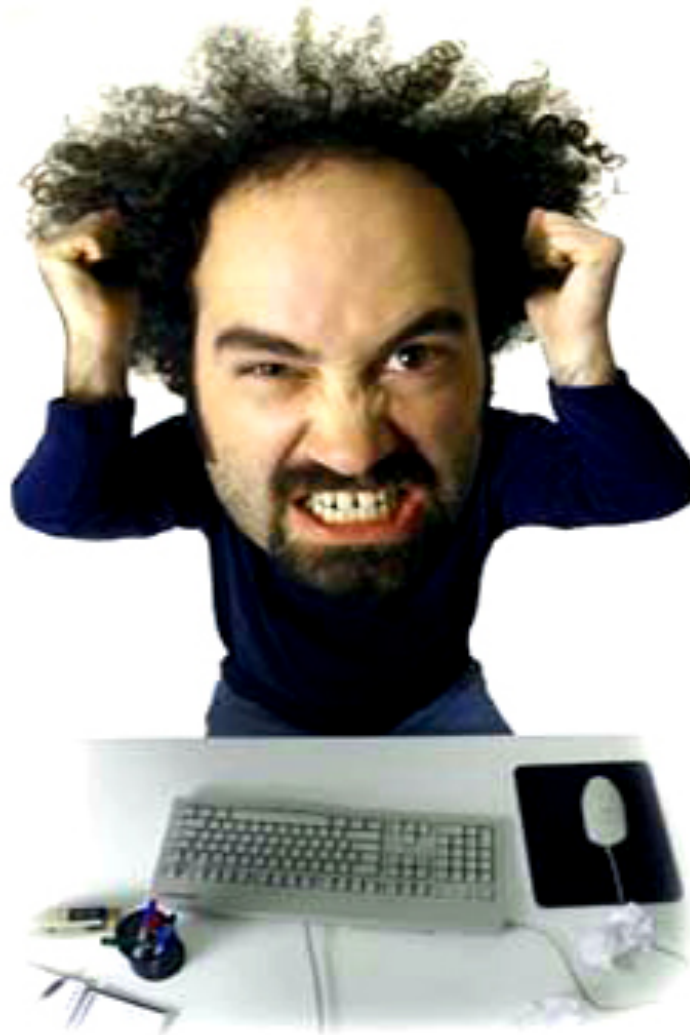
```
bin/plugin --url file:///path/to/plugin --install plugin-name
```

- Plugin management:

```
curl -XGET 'http://localhost:9200/_nodes/plugins'
```

```
plugin --remove <pluginname>
```

When to write a plugin?



Writing your own plugin: Gotchas

- **Maintenance** – the deeper you go in the API the harder it is to keep it up to date
- **Versioning and installation** on (large) clusters
 - Though can be solved using puppet, docker et al
- **Auxiliary data** (like dictionaries etc)
- **Testing & Debugging**

Code: Writing your own plugin

- JAR file with bootstrap code:

```
public class ExamplePlugin extends AbstractPlugin {
    @Override public String name() {
        return "example-plugin";
    }

    @Override public String description() {
        return "Example Plugin Description";
    }

    @Override
    public Collection<Class<? extends Module>> modules() {
        Collection<Class<? extends Module>> modules = Lists.newArrayList();
        modules.add(ExampleRestModule.class);
        return modules;
    }
}
```

- Embed this as es-plugin.properties:

`plugin=org.elasticsearch.plugin.example.ExamplePlugin`

Thank you.
Questions?

Itamar Syn-Hershko

<http://code972.com>

@synhershko