

# Elastic Agent: Ingesting log files, how hard can it be?

Jan 20, 2026

Denis Rechkunov





# Denis Rechkunov

Principal Software Engineer



Elastic Observability



Beats

<https://rdner.de>

# **Logs are a foundation of Observability**

# Let's talk about ingestion of log files

What's the general idea?



## Read a log line

We need a component that reads files line by line.



## Map it to JSON

We need to somehow convert it to JSON.



Elasticsearch

## Send it to Elasticsearch

We need to use an Elasticsearch client to send this JSON to an index.

**Raise your hand if you think  
this is easy** 

```
while read l
do curl -H"Authorization: ApiKey BASE64_API_KEY" \
  -XPOST localhost:9200/my-index/_doc \
  -HContent-Type:application/json \
  -d"$l"
done < data.ndjson
```

**Nope, it's hard.**

**Wouldn't it be fun to design  
it together really quick?**



## Let's split the problem into these 3 parts:



### Read a log line

We need a component that reads files line by line.



### Map it to JSON

We need to somehow convert it to JSON.



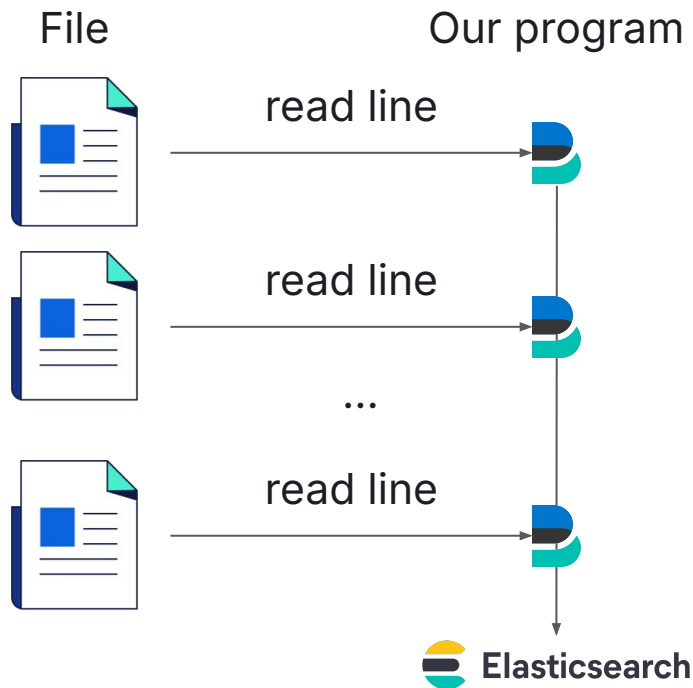
Elasticsearch

### Send it to Elasticsearch

We need to use an Elasticsearch client to send this JSON to an index.

# Part 1: Read a log line

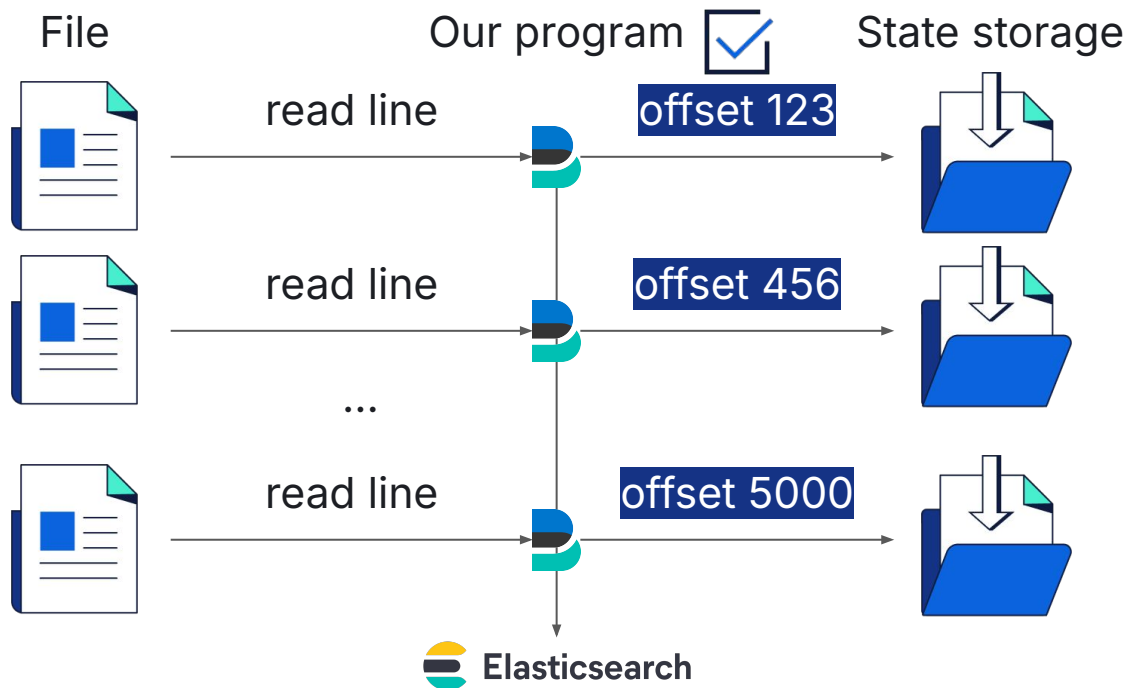
# Our First Approach





Credit: Heroes of Might and Magic III (videogame)

# Offset Tracking

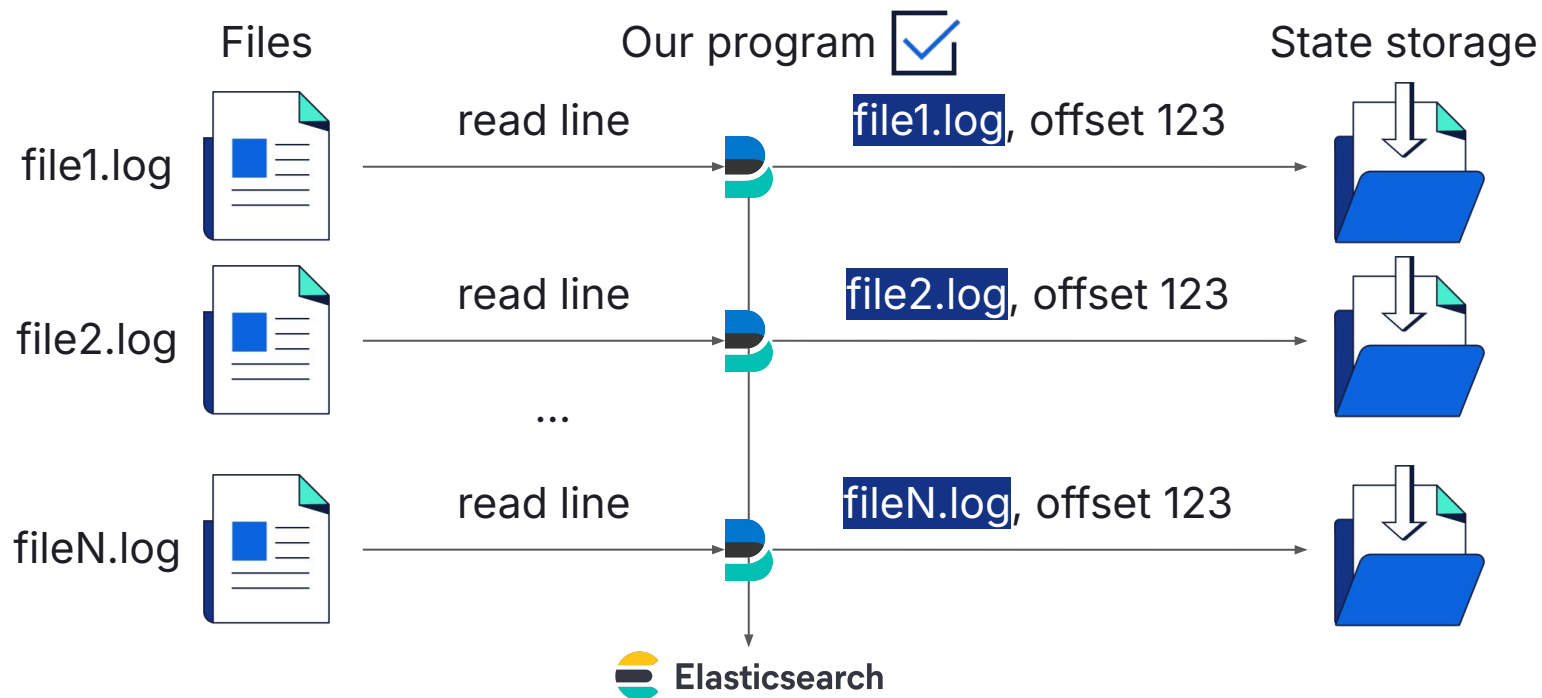






Credit: Heroes of Might and Magic III (videogame)

# Per File Offset Tracking



# Scaling Considerations

How to optimize for **thousands of files**

- State storage:
  - Key-Value disk storage.
  - Efficient writes, e.g. [Write-ahead logging](#) (WAL) with checkpointing.
  - In-memory state view or direct access to the on-disk storage if it's fast enough.
- Reading files:
  - OS has a file handle limit per process, sometimes as low as 256 (Mac).
  - Need to constantly close files on idle and re-open files on change.





Credit: Heroes of Might and Magic III (videogame)

# Log Rotation

Challenges to pick a **stable file identifier**

- Throughout the lifetime of the file we need to have a stable identifier.
- Files can be:
  - renamed/moved
  - truncated
  - compressed to Gzip
  - removed
- We cannot rely on its filename as a stable identifier.

# File Identification

Based on file system metadata

- It's kind of a **coordinate** where the file is, **not** its **name**.
- For **Unix**-like systems, it's the [stat](#) system call and a combination of:
  - device number (**st\_dev**)
  - [inode](#) (**ino\_t**)
- For **Windows**, it's the [file handler](#) and a combination of:
  - **dwVolumeSerialNumber**
  - **nFileIndexLow**
  - **nFileIndexHigh**

More details with examples at <https://www.elastic.co/docs/reference/beats/filebeat/file-identity>

# File Identification

Based on the content

- The content is what makes this file unique, nothing else affects the identity.
- Two main approaches:
  - **Hash** from the first **N bytes** – ingestion delay until the size is at least **N** bytes.
  - **Variable** length **match** until different – ingest right away, delay a file if collisions occur until it grows different.

More details with examples at <https://www.elastic.co/docs/reference/beats/filebeat/file-identity>

# File Identification

Why is file system metadata unreliable (Unix)?

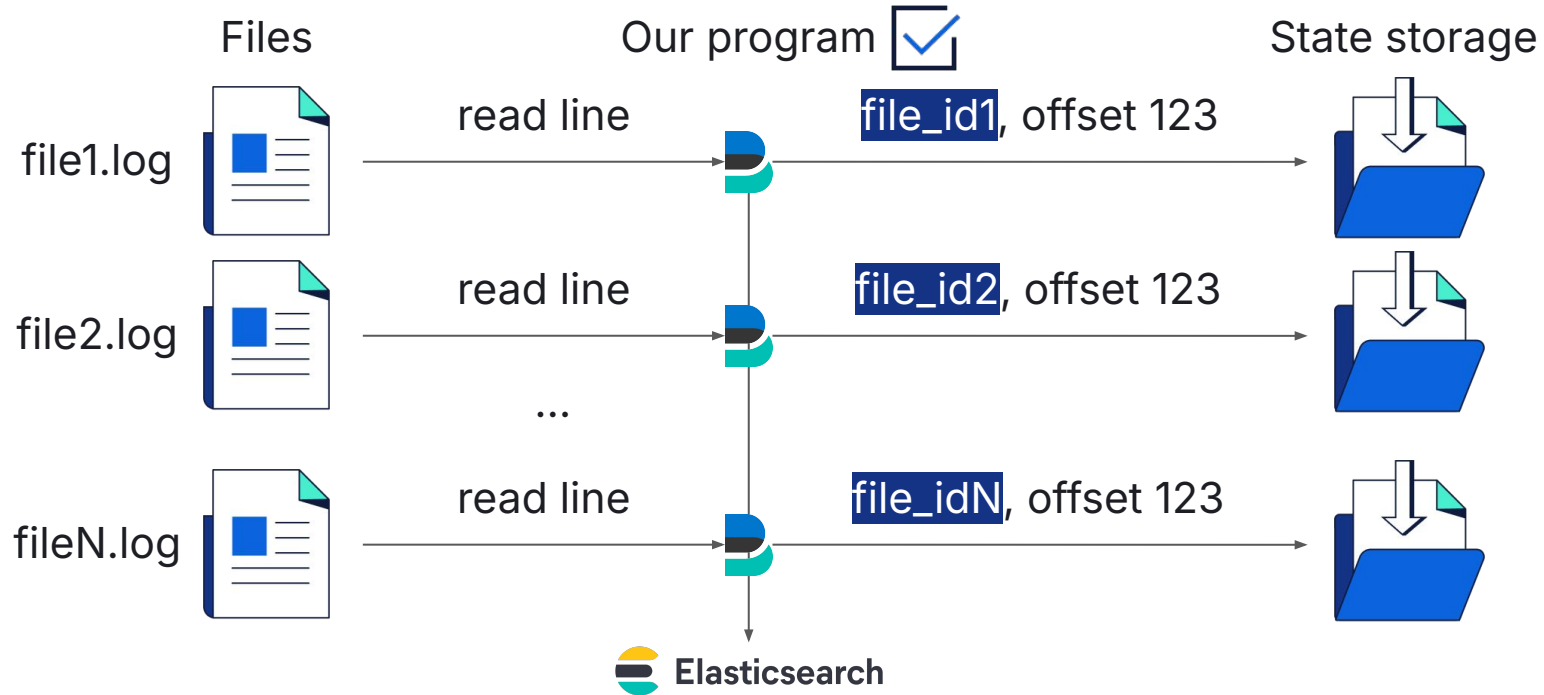
- Not every file system implements **inode** very well.  
*Examples:* NFS, FUSE filesystems.
- **inode** number gets reused for deleted files.  
*Examples:* containerized and virtualized environments.
- The **device ID/number** is unstable.  
*Examples:* when using the Linux [LVM](#) (Logical Volume Manager), device numbers are allocated dynamically at module load.  
(refer to [Persistent Device Numbers](#) in the Red Hat Enterprise Linux documentation)

# File Identification

Better use the content-based identifiers

- Files need to be different in content (most of log lines have timestamps, and it's enough).
- If the beginning of the file changed, it's fair to assume it's a new file.
- However, we have to delay ingestion of some files until we can create a **unique stable identifier** for them.

# Per File Offset Tracking with ID



## Part 2: Map it to JSON



# A log line to a JSON document

Straight-forward

- Some logs are already in JSON, nothing to do.
- Some logs are not JSON and they need to be converted  
e.g. with [Grok patterns](#) (not to be confused with Grok AI 🧑).
- 🎉



Credit: Heroes of Might and Magic III (videogame)

# Processors

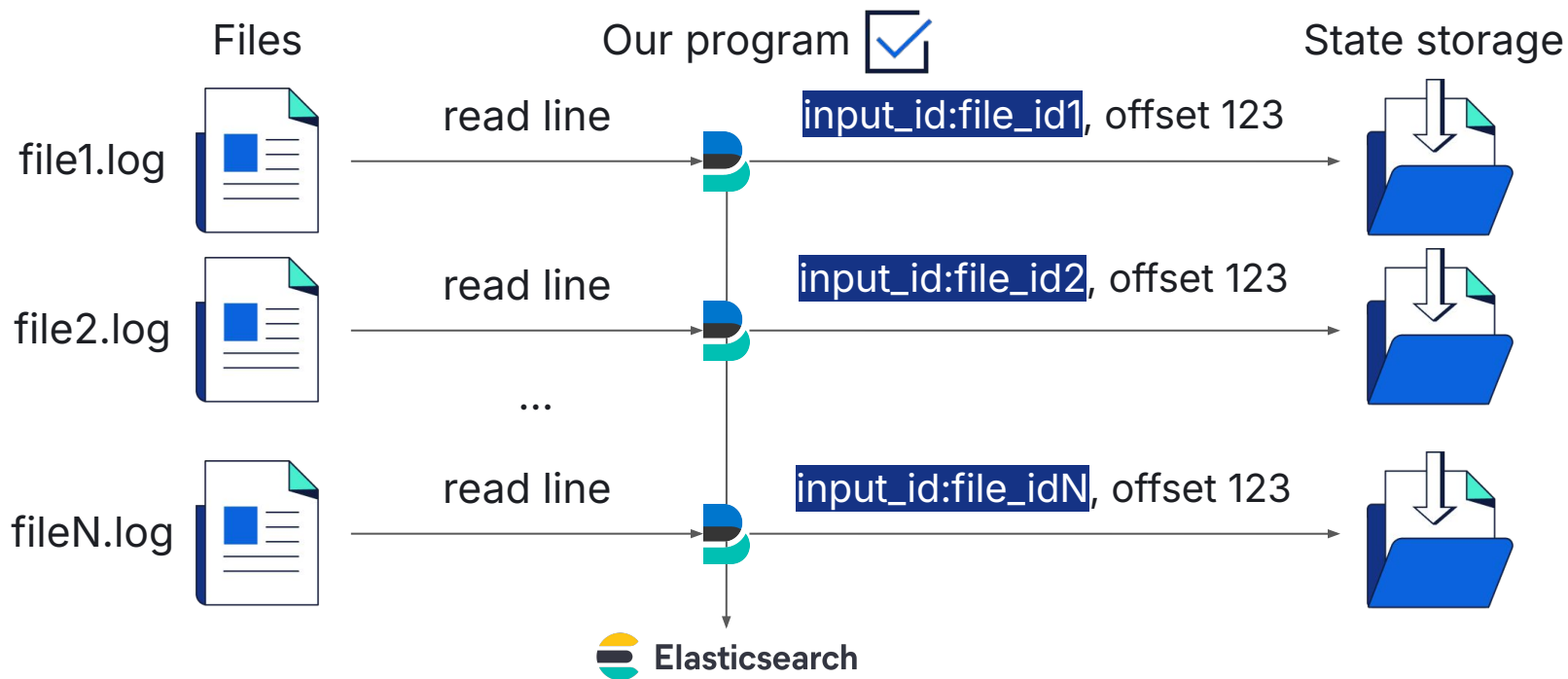
## Requirements

- We introduce a processing pipeline (aka **input**) with [dozens of processors](#).
  - add\_field
  - drop\_field
  - rename
  - decode
  - ... etc.
- Every processor must be:
  - lightweight – it's a **hot code path**
  - atomic – make backups before making changes
  - conditional – e.g. when/if/then/else



Credit: Heroes of Might and Magic III (videogame)

# Offset Tracking with File ID and Input ID



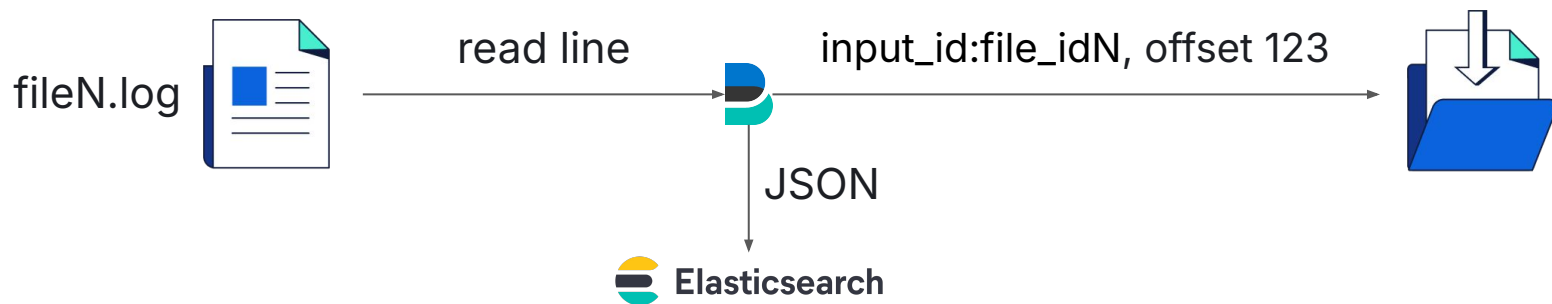


# **Part 3: Send it to Elasticsearch**

# Sending data to Elasticsearch

## Straight-forward

- Use the official client, e.g. for [Go](#).
- Send each JSON document to Elasticsearch.
- 🎉





Credit: Heroes of Might and Magic III (videogame)



# Sending to Elasticsearch

How to gain resilience

- Implement a queue: memory or disk.
- Combine JSON documents in batches and use the [bulk API](#).
- Track acknowledgements for documents.
- Store offset **after** a document is **acknowledged**.
- Do exponential back-offs and retries.



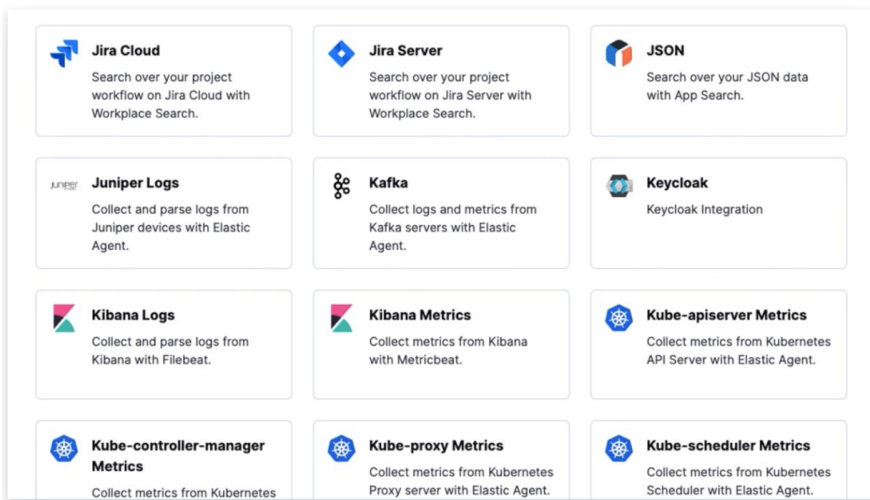
Credit: Heroes of Might and Magic III (videogame)

# Bonus Part: Integrations

# Integrations


Logs have established formats, we need to:

- Manage the catalog of integration packages
- Version them
- Configure them



More at <https://www.elastic.co/integrations>



The image shows the interface of the video game Heroes of Might and Magic III. A large, brown, parchment-like text box is centered on the screen, containing the text "Uh-oh... Users want centralized management for all of this across hundreds of machines." Below the text is a small icon of a yellow checkmark inside a dark rectangle. The background of the game is visible around the text box, showing a green landscape with a small pond, a red-roofed house, and a path. The game's UI elements are visible on the right and bottom. On the right, there is a minimap at the top, a unit and building selection panel in the middle, and a resource and status panel at the bottom. The bottom panel shows various resource icons and their counts: 14, 6, 24, 0, 0, 5, and 7430. It also shows the current month and week: "Month: 1, Week: 4".

**Uh-oh...**  
**Users want centralized  
management for all of this  
across hundreds of machines.**

Credit: Heroes of Might and Magic III (videogame)

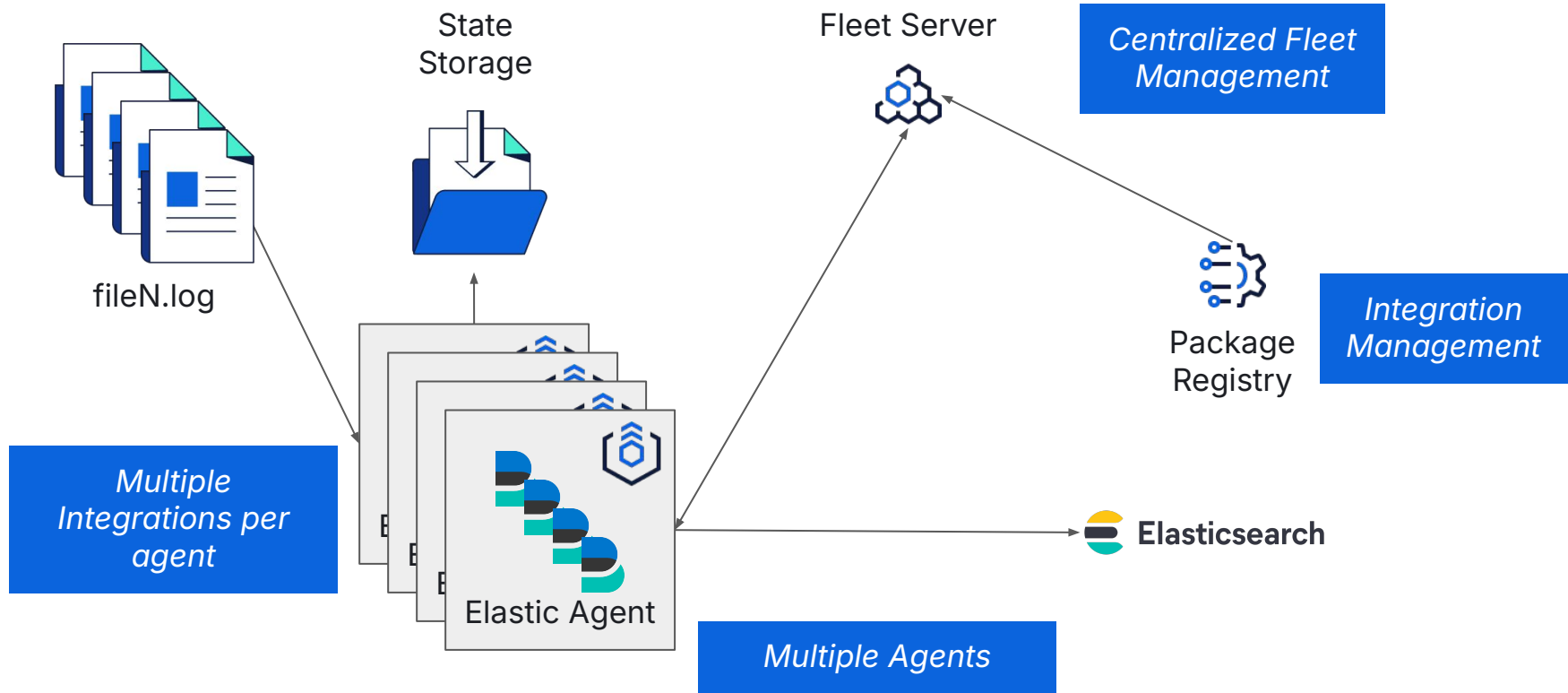
# Fleet Management

Centralized management across hundreds of machines

- Every ingestion program is configured and supervised by a locally running service – agent.
- A centralized server manages these agents.
- They can be re-configured with integration packages in runtime.
- Agents report their metrics and state back to the Fleet server.

More at <https://www.elastic.co/docs/reference/fleet>

# How do we solve it at Elastic?





Credit: Heroes of Might and Magic III (videogame)



# Elastic Agent as OpenTelemetry Collector

## Elastic Distribution of OpenTelemetry (EDOT)

- Elastic Agent is also an [OTel Collector](#).
- Beats (ingesting software running under the agent) are also [OTel receivers](#).
- OTel [components](#) can be running side by side with Elastic-specific components.
- Elastic is actively contributing to OpenTelemetry.
- Using open standards makes it future-proof and removes vendor-locking.



More at <https://www.elastic.co/docs/reference/edot-collector>

# Elastic Contributions to OpenTelemetry

As of January 2026:

OpenTelemetry Companies statistics (Contributions, Range: Last decade), bot... ▾		
Rank ^	Company	Number
	All	1301377
1	Splunk Inc.	290032
2	Microsoft Corporation	192429
3	LightStep Inc.	68899
4	Dynatrace LLC	63942
5	Elasticsearch Inc.	55090
6	Google LLC	46846
7	Amazon	43405
8	Datadog Inc	41711

Source <https://opentelemetry.devstats.cncf.io/d/5/companies-table>

**I guess the point is...**

IT'S DANGEROUS TO GO  
ALONE! TAKE THE AGENT.



Credit: Part of a series on The Legend of Zelda. (modified)

# Thank you!

My contacts



Slides

