# Mineração de Dados

## Data Sources & Data Collection (OSINT)

Miguel Rocha  | mrocha@di.uminho.pt

Diana Ferreira | diana@di.uminho.pt

# About me

## Education

**2016** Bachelor in Biomedical Engineering

**2019** MSc in Biomedical Engineering – Medical Informatics

**2024** PhD in Biomedical Engineering

## Experience

**2020 – 2025** Teaching at UMINHO (Databases, Knowledge Discovery, Neural Networks, Big Data)

**2019** Teaching Mobile Applications at IPCA

**2018** Member of the KEG at ALGORITMI Research Center

# Contents

# Learning Objectives

**By the end of this lecture, students should be able to:**

- Understand what OSINT is and where it applies;

- Know types of open data sources;

- Understand the OSINT cycle;

- Recognize practical applications and use cases

- Assess the quality, reliability, and legal risks of data;

- Perform basic OSINT data collection using open tools.

# Bibliography

**Books:**

1. "Open Source Intelligence Techniques" by Michael Bazzell

2. "The OSINT Handbook" by Dale Meredith

**Training Platforms:**

1. OSINT Framework (osintframework.com)

2. Trace Labs (for ethical practise)

3. Bellingcat's Online Investigation Toolkit

**Community and Resources:**

1. **Reddit:** r/OSINT

2. **GitHub:** Awesome OSINT (curated repository)

3. **Newsletters:** OSINT Weekly, Bellingcat

# The Data Age: Current Context

- 90% of the world's data was created in the last 2 years

- 328.77 million terabytes generated daily (2024)

- 5 billion internet users globally

- 70% of companies report that data is critical to decision-making

- Big Data Market: US$ 307.52 billion (2023)

" *In a world flooded with data, how can we find valuable information?*

# Why we collect data?

**Decision Making**

Evidence vs. Intuition

**Pattern Identification**

Trends and Correlations

**Problem Solving**

Accurate Diagnosis

**Optimization**

Process Improvement

**Innovation**

Insights for New Products/Services

**Competitiveness**

Strategic Advantage

# Data Quality Dimensions

**Accuracy:** The degree to which the data correctly represents the entity or atribute being described.
➤**How correct are the data values?**

**Completeness:** Amount of missing data from a given data set.
➤ **Is all information present?**

**Consistency:** Coherence in data between different sources/systems.
➤ **Does the data match other trusted sources?**

**Timeliness:** Data is up-to-date at the time of use.
➤ **How up-to-date is the data?**

**Validity:** Conformance to rules, formats, and standards.
➤ **Does the data conform to the predetermined format and constraints?**

**Uniqueness:** Ensures that there are no unnecessary duplications or overlappings within the data.
➤ **How much duplication is there in the records?**

# The Importance of Data Quality

**Decisions:** Bad data → Bad decisions

**Costs:** 20–35% of revenue is lost due to poor data quality

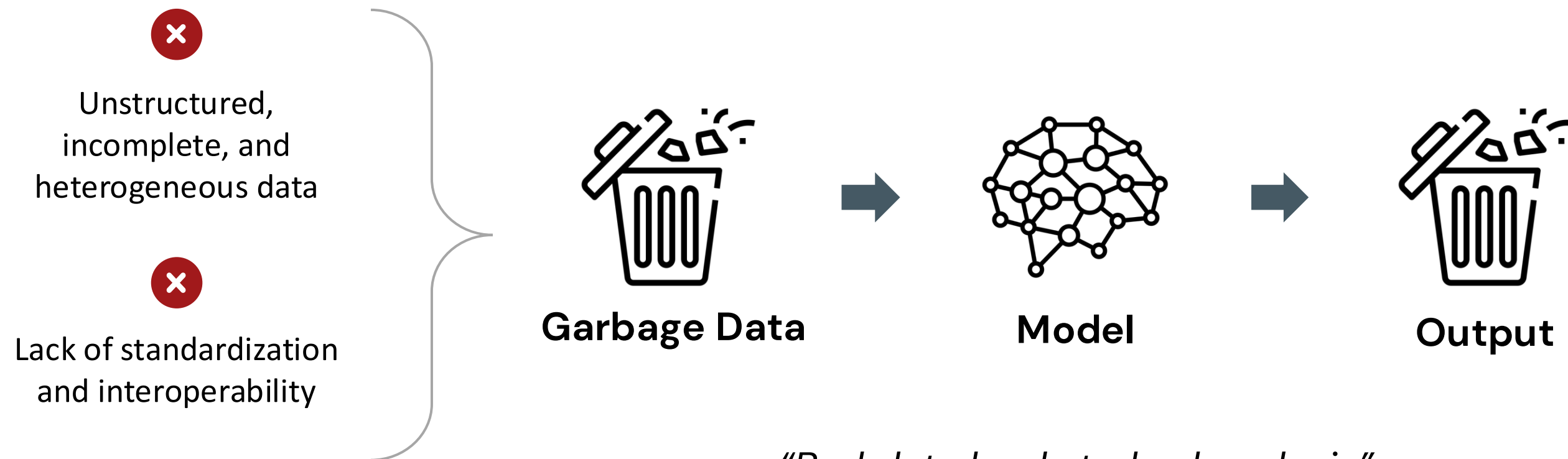**Efficiency:** 30% of analysts' time is spent cleaning data

**Compliance:** Fines for incorrect data

**Reputation:** Public errors cause damage to the brand

# GIGO Principle

## Garbage In, Garbage Out

❌ Unstructured, incomplete, and heterogeneous data

❌ Lack of standardization and interoperability

**Garbage Data** → **Model** → **Output**

*"Bad data leads to bad analysis"*

*"Biased data leads to biased conclusions"*

**EXAMPLE:** Incomplete sales data → Incorrect inventory forecasts → Loss of sales or excess inventory

💡 **Poor input = Worthless output**

# Modern Quality Challenges

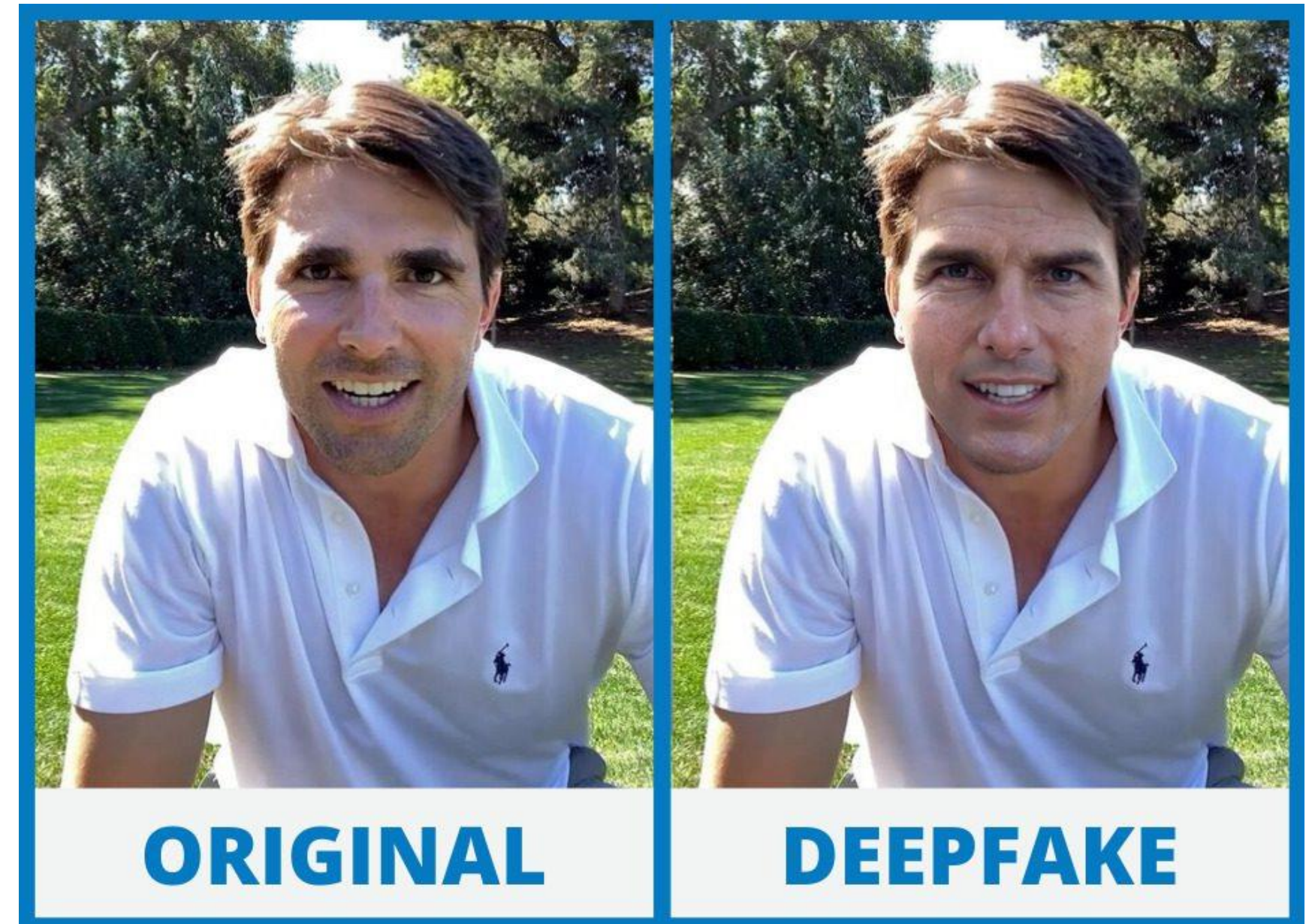## New Threats to Reliability

GENERATIVE AI AND DEEPFAKES:

– Convincing synthetic text, images, and audio

– Need for specialized detection tools

– Speed of creation vs. speed of verification

SYSTEMIC DISINFORMATION:

– Coordinated manipulation campaigns

– Closed ecosystems of disinformation

ENHANCED PRIVACY:

– Greater anonymization of personal data

– Difficulty in verifying identities



ORIGINAL    DEEPFAKE

# Data Sources vs Data Collection

**DATA SOURCES**

- **<u>Where</u>** data originates

- Passive existence

- These are the repositories, systems, platforms, or locations where data resides.

- They can be active (generating data) or passive (storing data).

<u>Examples of Sources:</u>
- Government databases
- Social media (Twitter, Facebook, LinkedIn)
- IoT sensors

**DATA COLLECTION**

- **<u>How</u>** data is obtained

- Active process

- It is the process/methodology of extracting data from sources

- It involves techniques, tools, and systematic approaches

<u>Examples of Collection Methods:</u>
Web scraping
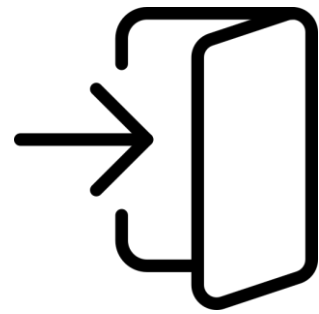Official APIs
Downloading open datasets

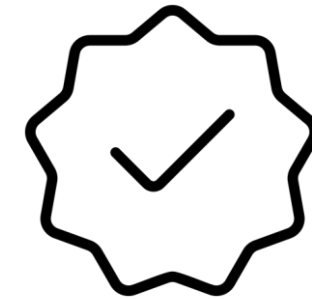⚠ **Poor source choice cannot be fixed by good collection**
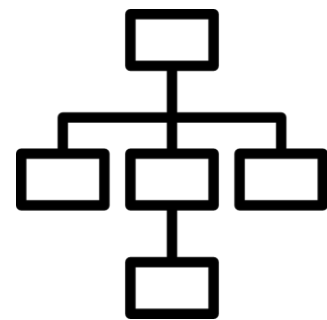
# Data Sources

Data sources can be classified by:

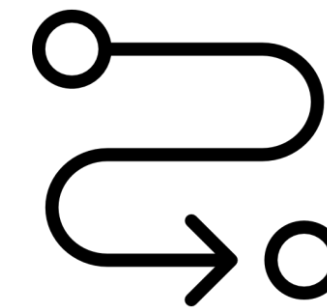**ACCESSIBILITY**
(open / semi-open / restricted)

**TEMPORALITY**
(static / dynamic)

**RELIABILITY**
(official / unofficial)

**STRUCTURE**
(structured / semi-structured / unstructured)

**ORIGIN**
(government / private / individual)

💡 **OSINT principle:** One source is no source ➜ Verification & Corroboration

# Data Sources

**Temporality**

**STATIC**

• Point in time

• Not automatically updated

• Specific versions

Examples:

Screenshots, Historical data files

**DYNAMIC**

• Constantly updated

• Continuous flow

• Real–time or near real–time

Examples:

Social media feeds, IoT sensor data

❌ **COMMON CHALLENGES:**

**Static:** Rapid obsolescence, conflicting versions
**Dynamic:** Data overload, difficulty storing historical data

# Data Sources

## Structure

|  |  |  |
|---|---|---|
| **STRUCTURED** | **SEMI-STRUCTURED** | **UNSTRUCTURED** |
| • Government open data | • APIs | • Web pages |
| • Statistical databases | • Logs | • News articles |
| • Corporate registries | • JSON / XML feeds | • Social media |
|  |  | • PDFs, images, videos |

# Data Sources

## Reliability

**OFFICIAL**

• Officially recognized entities

• Established verification processes

• Institutional accountability

Confidence Level:

HIGH (but not foolproof)

**UNOFFICIAL**

• No formal recognition

• Variable verification processes

• Diffuse responsibility

Confidence Level:

VARIABLE (requires intensive verification)

# Data Sources

**Accessibility**



### OPEN SOURCES (OSINT)
- Publicly available
- No significant access restrictions

Example: Government websites, public news

### SEMI–OPEN SOURCES
- Conditional access (login, registration)
- May have specific terms of use

Example: Social networks, registered forums

### CLOSED SOURCES
- Restricted or proprietary access
- Requires special authorization

Example: Private databases, internal systems

# Data Sources

**Origin**

### GOVERNMENT SOURCE

- Public funding
- Institutional accountability
- Transparency (in principle)

<u>Strengths</u>: Broad coverage, standardization

<u>Weaknesses</u>: Bureaucracy, slow updates

### PRIVATE (CORPORATE) SOURCE

- Private financing
- Focus on profit/value
- Intellectual property

<u>Strengths:</u> Innovation, speed, specialization

<u>Weaknesses:</u> Commercial bias, restricted access

### INDIVIDUAL SOURCE (CROWDSOURCED)

- Decentralized production
- Diversity of perspectives
- Variable quality

<u>Strengths</u>: Diversity, scale, cost

<u>Weaknesses</u>: Inconsistency, difficult to verify

# Data Sources

## Practical Evaluation Cases

|  | API SOCIAL MEDIA | IMAGES SOCIAL MEDIA |
|---|---|---|
| **ACCESSIBILITY** | Restricted (payment/authorization) | Semi-open (visual public, closed API) |
| **TEMPORALITY** | Dynamic | Static image, dynamic metadata |
| **RELIABILITY** | Official delivery system | Unofficial user content |
| **STRUCTURE** | Semi-structured (JSON) | Unstructured image + metadata |
| **ORIGIN** | Private corporate API | Individual user content |

# Data Collection

**Appropriate Strategies for Different Sources**

FOR WEB SOURCES:

Web scraping

REST/GraphQL APIs

RSS feeds

Change monitoring

FOR SOCIAL SOURCES:

Platform APIs (Twitter, Facebook)

Monitoring tools (Hootsuite, Brandwatch)

Manual scraping (when API is limited)

FOR DOCUMENTARY SOURCES:

OCR

Metadata extraction

Format conversion (PDF to text)

# Open Source Intelligence (OSINT)

## What is OSINT?

OSINT stands for Open–Source Intelligence, which refers to the collection of data from publicly available sources and performing intelligence techniques to gain useful and meaningful insights from it using a systematic methodology. These sources include social media platforms, news articles, government publications, academic papers, forums, blogs, and publicly accessible databases.

💡 **Key aspects:**

- Open ≠ free of constraints
- Intelligence ≠ raw data
- Methodology matters more than tools

☑ **OSINT adds:**

- Context
- Interpretation
- Judgment

# Open Source Intelligence (OSINT)

## What is OSINT?

| Concept | Focus |
|---|---|
| Open Data | Availability |
| Web Scraping | Collection technique |
| Data Mining | Pattern discovery |
| OSINT | Goal-oriented intelligence |
| Business Intelligence | Organizational decision-making |

**OSINT often feeds:**

• Web scraping

• NLP / text mining

• Graph databases

• Knowledge graphs
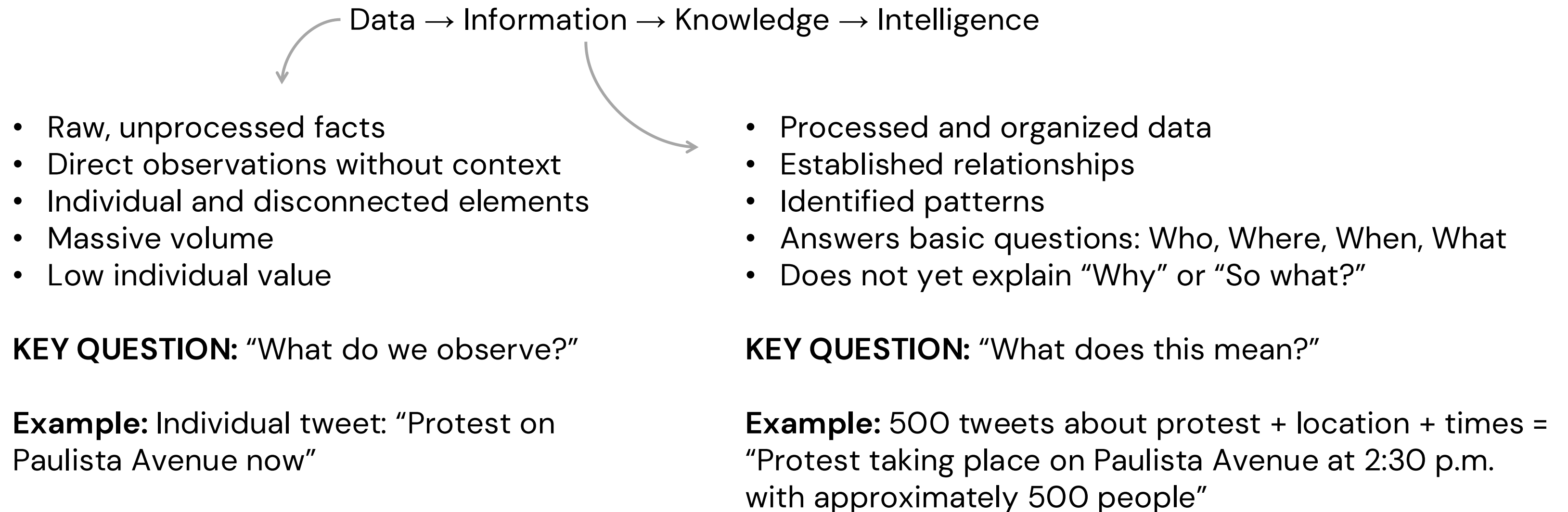
# Open Source Intelligence (OSINT)

**What is not OSINT?**

❌ It is not **hacking** (it does not involve unauthorized access)

❌ It is not **illegal surveillance** (it only uses public information)

❌ It is not just "**googling**" (it is methodical and systematic)

❌ It is not always **free** (some sources require payment)

❌ It is not **foolproof** (it requires verification and confirmation)

# Open Source Intelligence (OSINT)

## From Data to Intelligence

Data → Information → Knowledge → Intelligence

- Raw, unprocessed facts
- Direct observations without context
- Individual and disconnected elements
- Massive volume
- Low individual value

**KEY QUESTION:** "What do we observe?"

**Example:** Individual tweet: "Protest on Paulista Avenue now"

- Processed and organized data
- Established relationships
- Identified patterns
- Answers basic questions: Who, Where, When, What
- Does not yet explain "Why" or "So what?"

**KEY QUESTION:** "What does this mean?"

**Example:** 500 tweets about protest + location + times = "Protest taking place on Paulista Avenue at 2:30 p.m. with approximately 500 people"

# Open Source Intelligence (OSINT)

**From Data to Intelligence**

Data → Information → Knowledge → Intelligence

- Information analyzed and interpreted
- Implications and deeper meaning
- Basis for decisions
- Answers "Why?" and "So what?"
- Predictive or explanatory
- Action-oriented

**KEY QUESTION:** "What should we do with this?"

**Example:** Protest on Paulista Avenue + history of demonstrations + profile of protesters + content of posters = "Protest organized by unions against labor reform, with a tendency to grow and the possibility of clashes with the police"

- Application of intelligence with experience
- Understanding broader principles
- Strategic decision-making
- Includes values, ethics, experience
- Considers long-term consequences

**KEY QUESTION:** "What is the best course of action considering all factors?"

**Example:** Intelligence on protests + historical knowledge + organizational values + political context = "Do not intervene in the protest now, monitor and document; intervene only if there is violence."
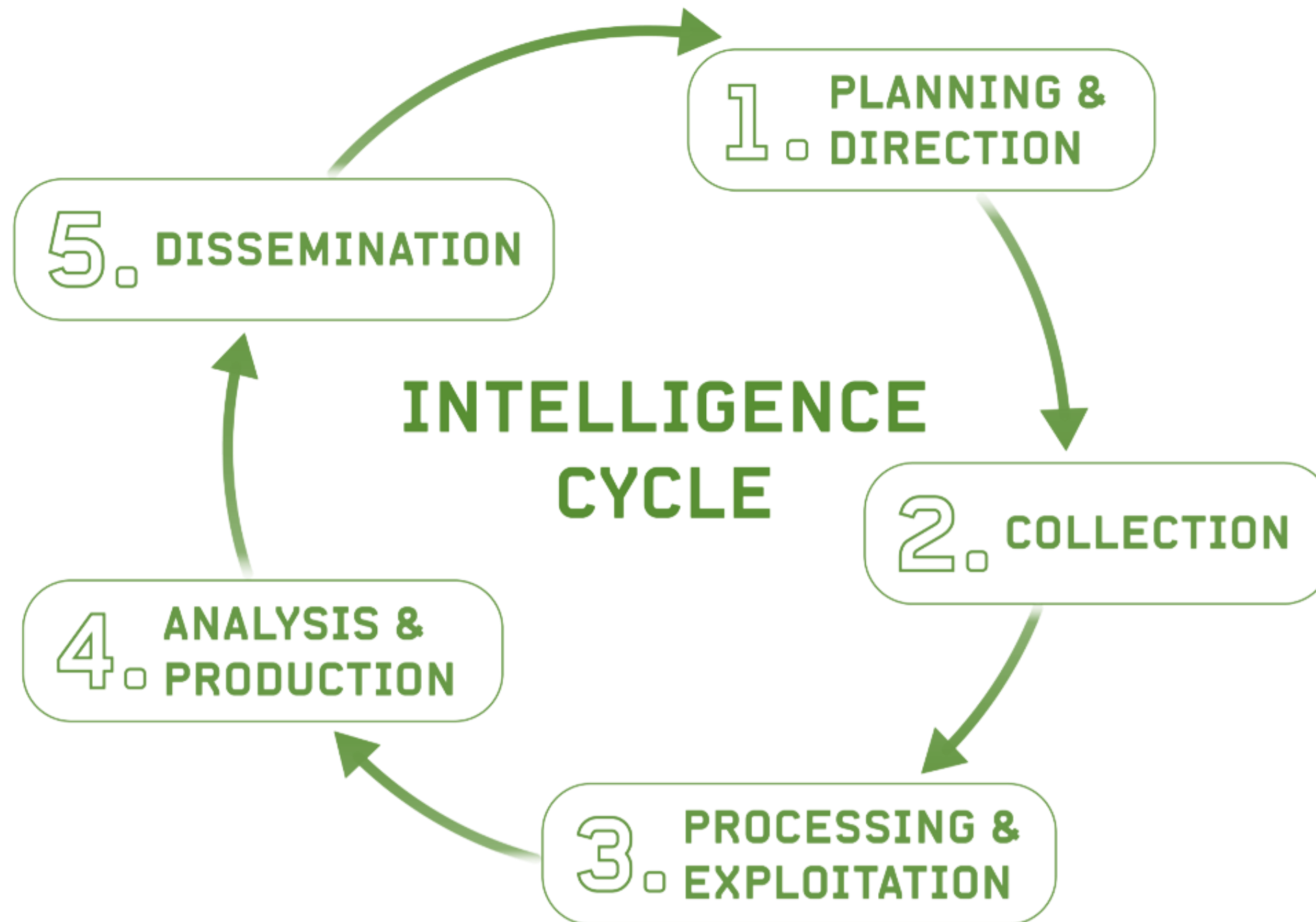
# Real World Use Cases

We can use OSINT almost everywhere, either its finding some online product for yourself or doing some bussiness anlysis for any corporate work.

- Investigating potentital security threats and vulnerabilities
- Monitoring and analyzing public opinions and sentiments
- Helping goverment agencices and military organizations in tracking terrirost activities
- Conducting market research and competitive analysis
- Identifying potential targets for sales and other bussiness development
- Helping law enformancements in missing person cases

# The OSINT Cycle

There is not a consensus  about the stages of the OSINT cycle, but the most common is:

**INTELLIGENCE CYCLE**

1. PLANNING & DIRECTION
2. COLLECTION
3. PROCESSING & EXPLOITATION
4. ANALYSIS & PRODUCTION
5. DISSEMINATION

⚠ Iterative, not linear

⚠ Collection without planning = noise

# The OSINT Cycle

➡ **<u>Planning & Direction</u>**

- Define intelligence requirements
- Identify questions to be answered
- Establish scope and boundaries
- Determine available resources
- Develop an operational plan
- Consider ethical and legal aspects

Apply the "Five W's and One H" framework to guide planning:

- **Who:** Identify who will conduct the work and who is the target of collection.
- **What:** Define the type, format, and content of data to be collected.
- **When:** Establish the timeframe for data collection and operational timing to avoid detection.
- **Where:** Identify relevant sources (e.g., news, social media, deep web, dark web) to investigate.
- **Why:** Link the information directly to mission objectives.
- **How:** Plan for security, access control, data storage, tools to be used, and archiving.

# The OSINT Cycle

➡ **Collection**

- Methodical: Systematic and organized
- Documented: Traceability of sources
- Ethical: Respect for terms of service
- Legal: Compliance with regulations
- Efficient: Optimized use of resources
- Comprehensive: Multiple sources and perspectives

**Often the longest step of the cycle, as many analysts rely on manual methods to conduct collection.**

# The OSINT Cycle

➡ **Processing**

- Organization of collected data

- Cleaning and standardization

- Metadata extraction

- Categorization and indexing

- Preparation for analysis

- Authenticity verification

# The OSINT Cycle

➡️ **Analysis**

- Link analysis

- Timeline analysis

- Spatial/geolocation analysis

- Social network analysis (SNA)

- Comparative analysis

- Pattern analysis

# The OSINT Cycle

## ➡️ Dissemination

- Target audience (who will receive the information)

- Appropriate format (report, dashboard, briefing)

- Level of detail required

- Frequency

- Distribution channels

- Protection of sources and methods

# Open Source Intelligence (OSINT)

The phrase "open source" does not refer to the open–source software movement, though many OSINT technologies do; rather, it refers to the open nature of the data being analyzed. Open source refers to the data/information that is pubicly available, it does not matter if its availiable in its online or offline mode.

**Examples of online information sources:**

– Search engines like Google, Yahoo, Bing

– Social media websites like Facebook, Instagram, Twitter, Reddit

– Sharing and publishing websites like Youtube, Pinterest, Medium

**Examples of offline information sources:**

– Goverment and law enforcement records

– Academic research and journals

– Annual Reports, Press Conferences

– Mass Media like TV, Newspaper

# OSINT Tools

**A. Search Engines & Discovery**

- **Google Dorking:** Using advanced operators (site:, filetype:, intitle:, inurl:) for deep web searches.

- **Shodan:** Search engine for IoT devices (servers, cameras, etc.).

- **Censys:** Similar to Shodan; indexes hosts and certificates (SSL certificates, hosts, etc.)

- **Searchcode:** Searches for source code across repositories.

- **Wayback Machine ([archive.org](archive.org)):** Views historical snapshots of websites.

- **Social-searcher.com:** Aggregated search

- **TweetDeck:** Real-time monitoring

- **Followerwonk:** Twitter follower analysis

# OSINT Tools

**B. Reconnaissance & Enumeration**

- **Maltego:** Powerful data mining and link analysis tool for visualizing relationships.
- **theHarvester:** Gathers emails, subdomains, hosts, and employee names from public sources.
- **Recon-ng:** Full-featured web reconnaissance framework.
- **SpiderFoot:** Automates OSINT collection from over 100 data sources.
- **OSINT Framework (osintframework.com):** A web-based directory of tools organized by category.

# OSINT Tools

**C. Social Media Intelligence (SOCMINT)**

- **Sherlock:** Searches for usernames across hundreds of social media sites.

- **Social Links:** A professional suite with both automated and manual search capabilities.

- **TweetDeck & Advanced Twitter Search:** For monitoring Twitter/X.

- **ImportYeti:** Tracks supply chains by searching global import/export records (useful for company research).

# OSINT Tools

**D. Domain & IP Investigation**

- **WHOIS Lookup** (e.g., ICANN Lookup, whois command): Finds domain registration details.

- **DNS Dumpster:** Discovers subdomains, DNS records, and related hosts.

- **VirusTotal:** Analyzes suspicious files, URLs, domains, and IPs for malware; also shows passive DNS data.

- **AbuseIPDB:** Checks IP addresses against a global blacklist for malicious activity.

- **URLScan.io:** Scans and analyzes websites, providing screenshots, tech stack, and associated links.

# OSINT Tools

**E. Geospatial Intelligence (GEOINT) & Imagery**

- **Google Earth Pro:** For historical imagery and advanced measurements.

- **Satellite Imagery:** Tools like **Google Maps**, **Bing Maps**, and specialized services like **Maxar**.

- **Flight Tracking: FlightRadar24** (live), **ADS-B Exchange** (unfiltered).

- **Marine Traffic:** For tracking ships and vessels.

# OSINT Tools

**F. Data & Document Analysis**

- **ExifTool:** Reads, writes, and edits metadata in files (images, PDFs, etc.).

- **PDF Metadata Analyzers:** Tools to inspect document properties and hidden data.

- **Have I Been Pwned (HIBP):** Checks if emails or passwords have been compromised in data breaches.

# OSINT Tools

**G. People & Identity Research**

- **Pipl:** One of the most powerful people search engines.

- **TruePeopleSearch:** US-focused people and phone number lookup.

- **FamilyTreeNow:** Genealogy-based people search.

- **LinkedIn:** A primary source for professional background and connections.

# OSINT Tools

**H. Network & Scanning Tools**

- **Nmap:** The standard for network discovery and security auditing.

- **Wireshark:** Network protocol analyzer for deep packet inspection.

- **FOCA:** Fingerprinting organizations with collected archives (analyzes metadata from documents).

# OSINT Tools

**I. Automation & Data Aggregation Platforms**

- **SpiderFoot:** (Also in Recon) Automates data collection from a vast array of sources.

- **IntelTechniques Tools:** A suite of online search tools created by Michael Bazzell for comprehensive searches.

- **Hunchly:** An OSINT capture tool that automatically documents and archives every webpage visited during an investigation.

# OSINT Tools

**Important Considerations & Best Practices:**

- **Legality & Ethics:** Always use these tools within legal and ethical boundaries. Respect Terms of Service and privacy laws.

- **Source Verification:** Information found via OSINT is not always accurate. Corroborate findings with multiple sources.

- **Operational Security (OPSEC):** Use VPNs, virtual machines, and separate accounts to protect your identity and avoid contaminating your investigation.

- **Tool Selection:** The right tool depends entirely on your specific requirement (e.g., finding a person, mapping a network, tracking social media).

- **Skill Over Tool:** A tool is only as effective as the operator. Understanding methodology and analytical thinking is more important than the tool itself.

# Data Quality in OSINT

## From Volume to Trusted Value

The Modern OSINT dilemma:

*"We have access to more data than ever before, but less certainty about its reliability."*

**Questions to ask for source reliability and bias:**

- Who produced this data?

- For what purpose?

- What is missing?

- Who benefits from this narrative?

⚠️ **Risks of Low-Quality Data:**

🚫 Decisions based on misinformation

🚫 Reputational damage from sharing misinformation

🚫 Time wasted on subsequent verification

🚫 Compromised operations/investigations

🚫 Legal consequences for using false information

# Data Quality in OSINT

## Trust Checklist

### 1. ORIGIN AND SOURCE

- ☑ Identifiable and traceable source
- ☑ History of source reliability
- ☑ Known motivations and biases
- ☑ Official channels vs. third parties

### 2. CORROBORATION

- ☑ Confirmed by multiple independent sources
- ☑ Temporal and factual consistency
- ☑ Supporting technical evidence
- ☑ Absence of logical contradictions

### 3. CONTEXTUALIZATION

- ☑ Environment and circumstances understood
- ☑ Cultural/regional elements considered
- ☑ Coherent timeline

# Data Quality in OSINT

**Reliability Scale**

LEVEL 1: HIGH RELIABILITY

📊 Official government data

🎓 Peer-reviewed academic research

🏛️ Public court documents

📰 Established news agencies with verification

# Data Quality in OSINT

**Reliability Scale**

LEVEL 2: MEDIUM RELIABILITY

📱 Verified social media profiles

🏢 Official corporate websites

📈 Reports from recognized companies

🗺️ Geospatial data from established sources

# Data Quality in OSINT

**Reliability Scale**

LEVEL 3: LOW RELIABILITY

👥 Anonymous social media accounts

💬 Unmoderated forums

🎭 Viral content with no clear source

⛔ Sources with a history of misinformation

# Data Quality in OSINT

**Framework S.A.F.E.**

## S
### SOURCE
- Identifiable origin?
- History of accuracy?
- Known biases?

## A
### ACCURACY
- Verifiable facts?
- Internal consistency?
- External corroboration?

## F
### FRESHNESS
- When was it created?
- When was it published?
- Is it still relevant?

## E
### EVIDENCE
- Technical support?
- Complete context?
- Evidence?

# Data Quality in OSINT

**Tools for Quality Assurance**

CONTENT VERIFICATION:

InVID / WeVerify (video analysis)

FotoForensics (image analysis)

RevEye (multi-platform reverse search)

SOURCE TRIANGULATION:

TinEye + Google Images + Yandex

BuzzSumo (virality analysis)

CrowdTangle (social media tracking)

CHANGE MONITORING:

Wayback Machine (web archive)

Versionista (website change monitor)

Changedetection.io (change alerts)

# Data Quality in OSINT

## Documentation Principles

For each piece of information collected:

- DATE/TIME of collection

- FULL URL with parameters

- METHOD of capture (print, API, download)

- TOOLS used

- CONTEXT of discovery

- Initial HYPOTHESES

- CONFIDENCE LEVEL assigned (1–5)

Example of label:

```
Source: Twitter/@oficial_profile
Collection: 2024-03-15 14:30 UTC
Method: Print screen + HTML file
Trust: 3/5 (verified profile, but
information not corroborated)
ID: OSINT-2024-015-TW-001
```

# Data Quality in OSINT

✅ **QUALITY > QUANTITY**
Five reliable sources are worth more than 50 dubious ones.

✅ **VERIFICATION IS A PROCESS, NOT AN EVENT**
Reassess continuously as new information emerges.

✅ **TRANSPARENCY BUILDS CREDIBILITY**
Document sources, methods, and levels of confidence

✅ **CONFIRMATION BIAS IS ENEMY NO. 1**
Actively seek information that contradicts your hypotheses

✅ **TOOLS HELP, HUMANS DECIDE**
AI helps, but critical human analysis is irreplaceable

# Good Practices

☑ **Document everything:** Sources, methods, dates –> If it cannot be explained, it cannot be trusted

☑ **Always verify:** Confirm with multiple sources

☑ **Stay organized**: Use management systems

☑ **Protect your identity**: When appropriate

☑ **Respect ToS**: Platform terms of servisse

☑ **Constantly evolve**: New tools emerge

☑ **Collaborate**: The OSINT community shares knowledge

# Open ≠ Ethical ≠ Legal

☒ Key **misconception**:

*"If it's public, it's safe to use."*

Reality:

• Legal ≠ ethical

• Accessible ≠ harmless

• Data aggregation increases risk

# Legal Aspects

- Public personal data is still personal data

- Specific and legitimate purpose

- Data minimization (collect only what is necessary)

- Legal basis for processing

- Rights of data subjects (access, rectification, deletion)

- International impact (GDPR has extraterritorial reach)

# Ethical Aspects

- **Legality:** Never violate laws

- **Proportionality:** Means appropriate to the ends

- **Transparency:** Document methods and sources

- **Accountability:** Be accountable for actions

- **Respect:** For rights and privacy

- **Good ends:** Legitimate and socially accepted purposes

💡 Ethical OSINT is **defensible OSINT!**

# Limitations of OSINT

- Incomplete visibility

- Deception & misinformation

- Platform dependency

- Legal constraints

- Analyst bias

# Summary and Key Takeaways

- OSINT is a **methodology**, not a toolset

- Source selection is critical

- Ethics and legality are central

- Analysis creates value, not collection

- OSINT integrates naturally with data science

# Application to the KE Project

- **FOR THE DATA INPUT PHASE (Requirement: >100k records)**

**A) Financial Domain:**
1. Public finance APIs:
  – Yahoo Finance API (free tier)
  – Alpha Vantage (free: 5 calls/min)
  – Twelve Data (free tier)

2. News sources:
  – NewsAPI.org (free: 100 requests/day)
  – RSS feeds from economic newspapers
  – Twitter API for financial trending topics

3. Government data:
  – Bank of Portugal (open data)
  – CMVM (Portuguese Securities Market Commission)
  – INE – Economic statistics

**B) Health Domain:**
1. Public datasets:
  – Kaggle medical datasets
  – UCI Machine Learning Repository
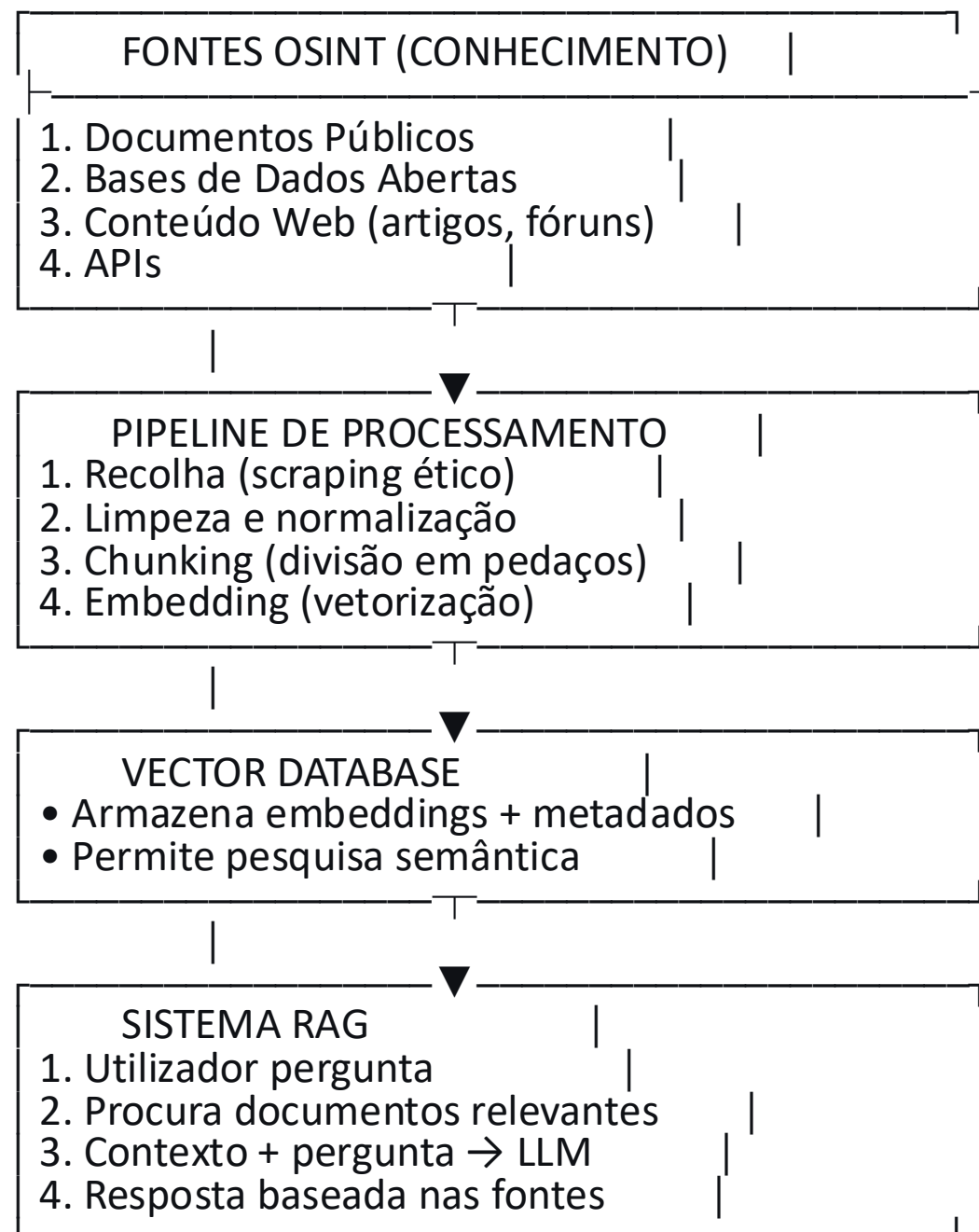  – Data.gov.pt (public health)

2. Medical literature:
  – PubMed API (free)
  – ClinicalTrials.gov API
  – Orphanet (rare diseases)

3. Simulated sensor data:
  – Generate synthetic data based on real patterns
  – Use public fitness APIs (Fitbit, Google Fit – limited)

# Application to the KE Project

- **FOR THE RAG SYSTEM (Retrieval–Augmented Generation)**

```
┌─────────────────────────────────────┐
│      FONTES OSINT (CONHECIMENTO)    │
├─────────────────────────────────────┤
│ 1. Documentos Públicos              │
│ 2. Bases de Dados Abertas           │
│ 3. Conteúdo Web (artigos, fóruns)   │
│ 4. APIs                             │
└──────────────────┬──────────────────┘
                   │
                   ▼
┌─────────────────────────────────────┐
│     PIPELINE DE PROCESSAMENTO       │
│ 1. Recolha (scraping ético)         │
│ 2. Limpeza e normalização           │
│ 3. Chunking (divisão em pedaços)    │
│ 4. Embedding (vetorização)          │
└──────────────────┬──────────────────┘
                   │
                   ▼
┌─────────────────────────────────────┐
│         VECTOR DATABASE             │
│ • Armazena embeddings + metadados   │
│ • Permite pesquisa semântica        │
└──────────────────┬──────────────────┘
                   │
                   ▼
┌─────────────────────────────────────┐
│          SISTEMA RAG                │
│ 1. Utilizador pergunta              │
│ 2. Procura documentos relevantes    │
│ 3. Contexto + pergunta → LLM        │
│ 4. Resposta baseada nas fontes      │
└─────────────────────────────────────┘
```

✅ **ADVANTAGES**

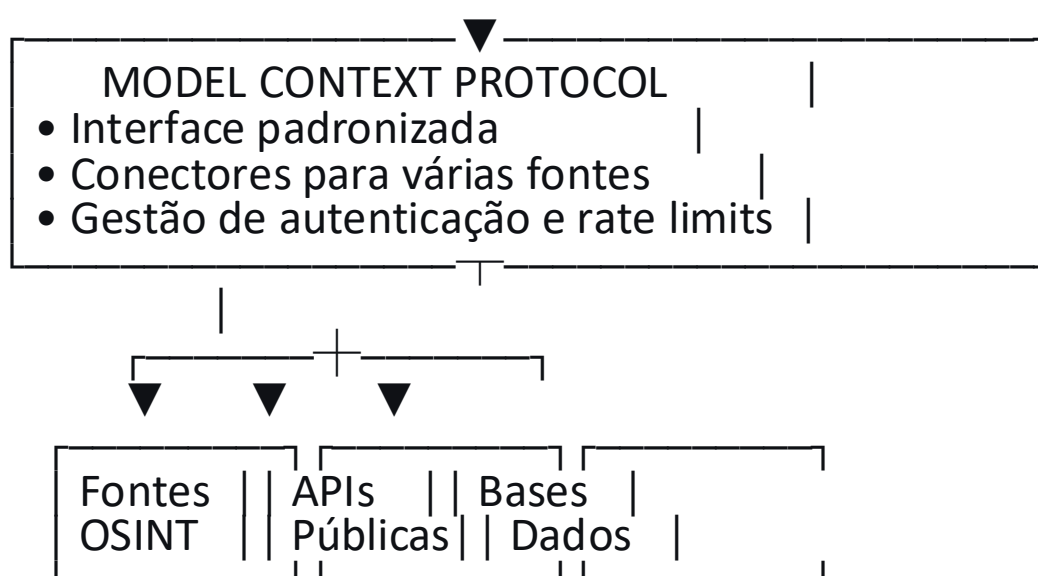**Up-to-date:** OSINT sources guarantee recent information

**Local context:** Specific data from Portugal

**Zero cost:** All sources are public

**Transparency:** Original sources can be cited

# Application to the KE Project

- **FOR THE MCP (MODEL CONTEXT PROTOCOL)**

```
                      ▼
   ┌──────────────────────────────────────┐
   │    MODEL CONTEXT PROTOCOL        │    │
   │  • Interface padronizada         │    │
   │  • Conectores para várias fontes │    │
   │  • Gestão de autenticação e rate limits │ │
   └──────────────────────────────────────┘
              │       ┬
              │       +
          ▼       ▼       ▼
   ┌──────┐┌──────┐┌──────┐
   │Fontes││ APIs ││ Bases│ │
   │OSINT ││Públicas││Dados│ │
   └──────┘└──────┘└──────┘
```

✅ **ADVANTAGES**

**Extensibility:** Easy to add new OSINT sources

**Modularity:** Each source is a separate "plugin"

**Resilience:** If one API fails, others continue to function

**Performance:** Cached data + real–time updates

# Practical Application

➡️ **Worksheet PL2**