

# Project Statement for Milestone 2

Kongqueror

Ross Kugler, Ben Bordon, Huy Ky

## Overview:

At the end of Milestone 2 of the project, teams should have prepared the data, including reducing, cleansing, and transforming, for storage and processing. Teams should also have chosen a NoSQL database, designed a non-relational schema, and ingested the (reduced) dataset into the database. The team should have validated the data ingestion process using appropriate database queries.

Teams should not use a hosted and managed NoSQL database service for this milestone.

Teams are not required to perform distributed data processing using Hadoop or Spark in Milestone 2. They are, however, expected to choose a NoSQL database that integrates well with Hadoop and/or Spark.

All team members are expected to make a significant contribution to the project milestone tasks.

Note: Since some datasets are significantly larger than others, teams are recommended to work with reduced dataset (>10 MB and <500 MB) during Milestones 1-3. Team would, however, still need to implement the final solution on a large dataset (>1 GB and <20 GB).

## Project Report Topics:

The report should cover the following subtopics and answer the questions listed:

1. Data Preparation:
  - a. Describe the data reduction, data cleansing and data transformation steps you have performed so far. Include pseudo-code for each of these steps in your description.

### Reduction:

Due to the nature of the dataset, we can just take the first few data crawl zip files. We will aim to keep the dataset <500MB for the first few milestones.

### Cleansing:

- Missing values:
  - Fill missing rates or ratings with 0
  - Drop rows with too many missing fields.
- Invalid values:
  - Check that age, length, views, ratings, comments are non-negative.
  - Remove or correct rows where length is unrealistically high (e.g., 999999).
- Data type enforcement:

- Convert columns to consistent types:
  - IDs → string
  - views/ratings/comments → integer
  - rate → float
- Normalization:
  - Standardize text categories (e.g., lowercase all category names, remove trailing spaces).

Transformation:

Raw data format:

video ID, uploader, age, category, length, views, rate, ratings, comments, related ID

To transform the raw data, turn it into a JSON-like document. The related ids need to be in an array.

-transform age into uploaded date: upload\_date = 2007-02-15 - age(days)  
 -put the related ids into an array in the video document

- b. You may have developed a parser to transform the raw data into semi-structured data. Briefly describe the parser algorithm with the help of a pseudo-code or source code snippet.

The algorithm processes each tab-separated record by first validating structural integrity (sufficient fields, valid YouTube video ID format) then applying data cleaning operations including safe type conversion with fallback defaults, value clamping for ratings (0-5 scale), category normalization, and relationship extraction from related video lists. The parser maintains comprehensive statistics throughout processing and employs defensive programming techniques with try-catch blocks and safe conversion methods to ensure robust handling of inconsistent raw data while preserving maximum data fidelity through minimal filtering approaches.

```
def parse_crawl_record(self, parts: List[str], crawl_id: str) -> Tuple[bool, Dict]:
    """Parse and clean a single crawl record"""
    # Check for malformed data
    is_malformed, reason = self.is_malformed_record(parts)
    if is_malformed:
        self.global_stats['malformed_removed'] += 1
        return False, {'error': f'malformed_{reason}'}

    video_id = parts[0].strip()
```

```

# Check for duplicates
if self.is_duplicate_record(video_id, crawl_id):
    self.global_stats['duplicates_removed'] += 1
    return False, {'error': 'duplicate'}


# Clean and structure the data
try:
    uploader = parts[1].strip() if parts[1] else "unknown"
    age_days = self.safe_int(parts[2])
    category = self.normalize_category(parts[3])
    length_sec = self.safe_int(parts[4])
    views = max(0, self.safe_int(parts[5])) # Ensure non-negative
    rate = max(0.0, min(5.0, self.safe_float(parts[6]))) # Clamp 0-5
    ratings = max(0, self.safe_int(parts[7]))
    comments = max(0, self.safe_int(parts[8]))

    # Clean related video IDs
    related_ids = []
    for rid in parts[9:]:
        if rid and self.validate_video_id(rid.strip()):
            related_ids.append(rid.strip())

    record = {
        'video_id': video_id,
        'uploader': uploader,
        'age_days': age_days,
        'upload_date': self.age_to_date(age_days),
        'category': category,
        'length_sec': length_sec,
        'views': views,
        'rate': rate,
        'ratings': ratings,
        'comments': comments,
        'related_ids': related_ids[:50], # Limit to prevent oversized
docs
        'crawl_id': crawl_id,
        'processed_at': datetime.now().isoformat()
    }

```

## 2. Database System:

- What NoSQL database are you using to store the data? Does the database system scale well with the data?

## MongoDB

Reasoning:

- Integrates with Apache Spark: MongoDB Spark Connector
    - The official integration maintained by MongoDB
    - Allows Spark to read from and write to MongoDB collections as if they were Spark DataFrames
    - Can run Spark SQL queries or MLlib pipelines directly on MongoDB data
  - Integrates with Hadoop: MongoDB Hadoop Connector
    - Allows Hadoop MapReduce jobs to read/write MongoDB collections
    - Can use MongoDB as both an input source and output sink for Hadoop jobs
  - Easy to manage stored data: MongoDB can store everything in one self-contained document
  - Query power: filtering, aggregation, indexing
  - Scalability: has built-in load balancing (sharding). Data can be stored across multiple servers, and automatically moved across servers to improve access times.
- b. Describe the non-relational schema you have implemented for the data. Why is the schema an appropriate one for your project?

Non-relational schema:

4 collections in MongoDB

videos: stores static attributes for each video

```
{  
  "_id": "video_id", // from the data's video id  
  "uploader": "...",  
  "category": "...",  
  "length_sec": 213,  
  "date_uploaded": ISODate("2007-01-01T00:00:00Z"), // calculated from age  
  "seen_in_crawls": ["0222", "0301", "0302"]  
}
```

video\_snapshots

```
{  
  "_id": { "video_id": "dQw4w9WgXcQ", "crawl_id": "0222" },  
  "video_id": "dQw4w9WgXcQ",  
  "crawl_id": "0222",  
  "age_days": 1234,           // as defined: days since YouTube establishment  
  "category": "Music",
```

```

    "length_sec": 213,           // integer
    "views": 1234567,
    "rate": 4.62,              // float
    "ratings": 9876,
    "comments": 5432
}

edges:
{
    "_id": ObjectId(),
    "crawl_id": "0222",
    "src": "dQw4w9WgXcQ",      // the video containing the "related IDs" list
    "dst": "oHg5SJYRHA0"       // one related video id
}

crawls: handy for per-crawl reporting/filters
{
    "_id": "0222",
    "date": ISODate("2007-03-01T00:00:00Z"),
    "notes": "crawled to fifth depth, but did not finish",
    "total_videos": 155513,
    "duration_sec": 93352
}

```

This schema is appropriate for our data because we want to analyze top-k, degree distribution, video category stats, and top video patterns. Storing videos allows for determining which videos are the most popular and their basic information. Storing edges makes the database much larger, but allows us to do graph-related queries and analysis with Spark GraphFrames. Storing the video snapshots allows us to see how videos evolve over time.

### 3. Data Ingestion and Query:

- a. Describe how you ingested the dataset into the database. Also, describe how you validated the ingestion step, perhaps through database queries.

Ingestion Process: data is processed and inserted into MongoDB using a python class that:

- Extracts ZIP files, reads crawl metadata, and processes data line-by-line.
- Filters out malformed records and duplicates.
- Converts each record into one of four MongoDB collections: videos, video\_snapshots, edges, and crawls.
- Inserts data in 1000-record batches using upserts and bulk operations.
- Skips previously processed crawls to maintain consistency.

Validation Methodology: Validation was done in three ways:

- **During ingestion** – tracked processing stats, checked document structure, and monitored batch success.
- **Post-ingestion** – used queries to compare document counts, check schema formats, assess data quality, and verify references.
- **Automated checks** – ran validate\_crawl\_processing.py to generate reports with counts, sample validations, quality metrics, and milestone compliance.

Provide performance results of your data ingestion and query operations. How would these results scale with data?

#### Performance Results:

```
2025-10-12 23:19:07,183 - INFO - [CONNECT] Connected to MongoDB:  
mongodb://localhost:27017/  
2025-10-12 23:19:07,183 - INFO - [MODE] Incremental processing mode - will skip  
already processed crawls  
2025-10-12 23:19:07,183 - INFO - [EXTRACT] Extracting crawl files...  
2025-10-12 23:19:07,184 - INFO - Extracting 0222.zip...  
2025-10-12 23:19:08,099 - INFO - Extracted 6 files to data\extracted\0222  
2025-10-12 23:19:08,100 - INFO - Extracting 0301.zip...  
2025-10-12 23:19:08,314 - INFO - Extracted 5 files to data\extracted\0301  
2025-10-12 23:19:08,314 - INFO - [SUCCESS] Extracted 2 crawl directories  
2025-10-12 23:19:08,319 - INFO - [INITIAL] No existing crawls found - processing  
all crawls  
2025-10-12 23:19:08,319 - INFO - [INDEX] Ensuring database indexes exist...  
2025-10-12 23:19:08,819 - INFO - [SUCCESS] Database indexes ensured  
2025-10-12 23:19:08,819 - INFO - [PROCESS] Processing crawl: 0222  
2025-10-12 23:19:08,820 - INFO - [LOG] Parsing log: log.txt  
2025-10-12 23:19:08,832 - INFO - Crawl ID: 0222  
2025-10-12 23:19:08,832 - INFO - Start: 0222 16:54:40  
2025-10-12 23:19:08,832 - INFO - Finish: 0227 17:04:29  
2025-10-12 23:19:08,836 - INFO - [ Processing 0.txt...  
Completed 0.txt: 189/189 valid (100.0%)  
2025-10-12 23:19:08,859 - INFO - [ Processing 1.txt...  
Completed 1.txt: 3,141/3,169 valid (99.1%)  
2025-10-12 23:19:10,211 - INFO - [ Processing 2.txt...  
Processed 10,000 records...  
2025-10-12 23:19:13,061 - INFO - Processed 20,000 records...  
2025-10-12 23:19:17,044 - INFO - Completed 2.txt: 23,454/23,543 valid  
(99.6%)  
2025-10-12 23:19:19,623 - INFO - [ Processing 3.txt...  
Processed 30,000 records...  
2025-10-12 23:19:21,195 - INFO - Processed 40,000 records...  
2025-10-12 23:19:25,494 - INFO - Processed 50,000 records...
```

2025-10-12 23:19:34,725 - INFO -	Processed 60,000 records...
2025-10-12 23:19:39,512 - INFO -	Processed 70,000 records...
2025-10-12 23:19:42,441 - INFO - (99.3%)	Progress: 50,000 lines, 49,653 valid
2025-10-12 23:19:44,367 - INFO -	Processed 80,000 records...
2025-10-12 23:19:49,029 - INFO -	Processed 90,000 records...
2025-10-12 23:19:53,575 - INFO -	Processed 100,000 records...
2025-10-12 23:19:58,346 - INFO -	Processed 110,000 records...
2025-10-12 23:20:03,209 - INFO -	Processed 120,000 records...
2025-10-12 23:20:05,651 - INFO - (99.0%)	Progress: 100,000 lines, 98,998 valid
2025-10-12 23:20:08,086 - INFO -	Processed 130,000 records...
2025-10-12 23:20:13,103 - INFO -	Processed 140,000 records...
2025-10-12 23:20:18,454 - INFO -	Processed 150,000 records...
2025-10-12 23:20:23,436 - INFO -	Processed 160,000 records...
2025-10-12 23:20:28,450 - INFO -	Processed 170,000 records...
2025-10-12 23:20:31,770 - INFO - (99.2%)	Progress: 150,000 lines, 148,762 valid
2025-10-12 23:20:34,623 - INFO -	Processed 180,000 records...
2025-10-12 23:20:40,174 - INFO -	Processed 190,000 records...
2025-10-12 23:20:45,630 - INFO -	Processed 200,000 records...
2025-10-12 23:20:51,009 - INFO -	Processed 210,000 records...
2025-10-12 23:20:56,570 - INFO -	Processed 220,000 records...
2025-10-12 23:20:59,742 - INFO - (99.2%)	Progress: 200,000 lines, 198,497 valid
2025-10-12 23:21:02,323 - INFO -	Processed 230,000 records...
2025-10-12 23:21:03,939 - INFO - (99.3%)	Completed 3.txt: 206,784/208,331 valid
2025-10-12 23:21:03,940 - INFO -	 Processing 4.txt...
2025-10-12 23:21:07,912 - INFO -	Processed 240,000 records...
2025-10-12 23:21:13,322 - INFO -	Processed 250,000 records...
2025-10-12 23:21:19,555 - INFO -	Processed 260,000 records...
2025-10-12 23:21:25,429 - INFO -	Processed 270,000 records...
2025-10-12 23:21:31,008 - INFO -	Processed 280,000 records...
2025-10-12 23:21:32,632 - INFO - (99.5%)	Progress: 50,000 lines, 49,753 valid
2025-10-12 23:21:37,056 - INFO -	Processed 290,000 records...
2025-10-12 23:21:42,859 - INFO -	Processed 300,000 records...
2025-10-12 23:21:48,443 - INFO -	Processed 310,000 records...
2025-10-12 23:21:54,203 - INFO -	Processed 320,000 records...
2025-10-12 23:22:00,188 - INFO -	Processed 330,000 records...
2025-10-12 23:22:01,306 - INFO - (99.4%)	Progress: 100,000 lines, 99,376 valid
2025-10-12 23:22:06,013 - INFO -	Processed 340,000 records...
2025-10-12 23:22:24,106 - INFO -	Processed 350,000 records...

2025-10-12 23:22:50,630 - INFO -	Processed 360,000 records...
2025-10-12 23:23:04,334 - INFO -	Processed 370,000 records...
2025-10-12 23:23:13,216 - INFO -	Processed 380,000 records...
2025-10-12 23:23:14,342 - INFO - (99.3%)	Progress: 150,000 lines, 148,968 valid
2025-10-12 23:23:31,722 - INFO -	Processed 390,000 records...
2025-10-12 23:23:38,344 - INFO -	Processed 400,000 records...
2025-10-12 23:23:47,466 - INFO -	Processed 410,000 records...
2025-10-12 23:23:55,489 - INFO -	Processed 420,000 records...
2025-10-12 23:24:29,944 - INFO -	Processed 430,000 records...
2025-10-12 23:24:34,791 - INFO - (99.3%)	Progress: 200,000 lines, 198,688 valid
2025-10-12 23:24:40,276 - INFO -	Processed 440,000 records...
2025-10-12 23:24:49,098 - INFO -	Processed 450,000 records...
2025-10-12 23:25:01,711 - INFO -	Processed 460,000 records...
2025-10-12 23:25:31,491 - INFO -	Processed 470,000 records...
2025-10-12 23:25:42,500 - INFO -	Processed 480,000 records...
2025-10-12 23:25:43,561 - INFO - (99.3%)	Progress: 250,000 lines, 248,366 valid
2025-10-12 23:25:50,265 - INFO -	Processed 490,000 records...
2025-10-12 23:26:11,714 - INFO -	Processed 500,000 records...
2025-10-12 23:26:46,340 - INFO -	Processed 510,000 records...
2025-10-12 23:26:54,751 - INFO -	Processed 520,000 records...
2025-10-12 23:27:02,891 - INFO -	Processed 530,000 records...
2025-10-12 23:27:03,713 - INFO - (99.3%)	Progress: 300,000 lines, 298,009 valid
2025-10-12 23:27:13,579 - INFO -	Processed 540,000 records...
2025-10-12 23:28:10,941 - INFO -	Processed 550,000 records...
2025-10-12 23:28:57,533 - INFO -	Processed 560,000 records...
2025-10-12 23:29:31,562 - INFO -	Processed 570,000 records...
2025-10-12 23:30:22,728 - INFO -	Processed 580,000 records...
2025-10-12 23:30:27,682 - INFO - (99.3%)	Progress: 350,000 lines, 347,628 valid
2025-10-12 23:31:37,335 - INFO -	Processed 590,000 records...
2025-10-12 23:32:03,843 - INFO -	Processed 600,000 records...
2025-10-12 23:32:32,004 - INFO -	Processed 610,000 records...
2025-10-12 23:33:09,297 - INFO -	Processed 620,000 records...
2025-10-12 23:33:27,460 - INFO -	Processed 630,000 records...
2025-10-12 23:33:28,734 - INFO - (99.3%)	Progress: 400,000 lines, 397,079 valid
2025-10-12 23:33:52,417 - INFO -	Processed 640,000 records...
2025-10-12 23:34:34,535 - INFO -	Processed 650,000 records...
2025-10-12 23:34:59,146 - INFO -	Processed 660,000 records...
2025-10-12 23:35:28,632 - INFO -	Processed 670,000 records...
2025-10-12 23:35:54,545 - INFO -	Processed 680,000 records...

2025-10-12 23:35:54,574 - INFO - (99.2%)	Progress: 450,000 lines, 446,447 valid
2025-10-12 23:36:31,600 - INFO -	Processed 690,000 records...
2025-10-12 23:36:47,474 - INFO -	Processed 700,000 records...
2025-10-12 23:37:02,125 - INFO -	Processed 710,000 records...
2025-10-12 23:37:10,228 - INFO -	Processed 720,000 records...
2025-10-12 23:37:18,001 - INFO - (99.2%)	Progress: 500,000 lines, 495,951 valid
2025-10-12 23:37:19,121 - INFO -	Processed 730,000 records...
2025-10-12 23:37:28,652 - INFO -	Processed 740,000 records...
2025-10-12 23:37:38,164 - INFO - (99.2%)	Completed 4.txt: 509,894/514,129 valid
2025-10-12 23:37:38,164 - INFO -	<b>■ Crawl 0222 summary:</b>
2025-10-12 23:37:38,164 - INFO -	Total lines: 749,361
2025-10-12 23:37:38,164 - INFO -	Valid records: 743,462
2025-10-12 23:37:38,165 - INFO -	Malformed removed: 5,899
2025-10-12 23:37:38,165 - INFO -	Duplicates removed: 0
2025-10-12 23:37:38,165 - INFO -	[SUCCESS] Completed crawl: 0222
2025-10-12 23:37:38,165 - INFO -	[PROCESS] Processing crawl: 0301
2025-10-12 23:37:38,166 - INFO -	[LOG] Parsing log: log.txt
2025-10-12 23:37:38,414 - INFO -	Crawl ID: 0301
2025-10-12 23:37:38,414 - INFO -	Start: 0228 23:53:41
2025-10-12 23:37:38,415 - INFO -	Finish: 0302 1:49:33
2025-10-12 23:37:38,417 - INFO -	<b>□ Processing 0.txt...</b>
2025-10-12 23:37:38,457 - INFO -	Completed 0.txt: 353/359 valid (98.3%)
2025-10-12 23:37:38,457 - INFO -	<b>□ Processing 1.txt...</b>
2025-10-12 23:37:41,311 - INFO -	Completed 1.txt: 3,614/3,636 valid (99.4%)
2025-10-12 23:37:41,311 - INFO -	<b>□ Processing 2.txt...</b>
2025-10-12 23:37:45,937 - INFO -	Processed 10,000 records...
2025-10-12 23:37:54,045 - INFO -	Processed 20,000 records...
2025-10-12 23:37:57,393 - INFO - (99.4%)	Completed 2.txt: 19,695/19,817 valid
2025-10-12 23:37:57,394 - INFO -	<b>□ Processing 3.txt...</b>
2025-10-12 23:38:02,260 - INFO -	Processed 30,000 records...
2025-10-12 23:38:17,248 - INFO -	Processed 40,000 records...
2025-10-12 23:38:24,839 - INFO -	Processed 50,000 records...
2025-10-12 23:38:32,780 - INFO -	Processed 60,000 records...
2025-10-12 23:38:40,565 - INFO -	Processed 70,000 records...
2025-10-12 23:38:43,075 - INFO - (99.7%)	Progress: 50,000 lines, 49,851 valid
2025-10-12 23:38:56,965 - INFO -	Processed 80,000 records...
2025-10-12 23:39:05,580 - INFO -	Processed 90,000 records...
2025-10-12 23:39:15,331 - INFO -	Processed 100,000 records...
2025-10-12 23:39:25,152 - INFO -	Processed 110,000 records...
2025-10-12 23:39:59,373 - INFO -	Processed 120,000 records...

```

2025-10-12 23:40:08,650 - INFO - Progress: 100,000 lines, 99,592 valid
(99.6%)
2025-10-12 23:40:14,833 - INFO - Processed 130,000 records...
2025-10-12 23:40:24,190 - INFO - Processed 140,000 records...
2025-10-12 23:40:45,828 - INFO - Processed 150,000 records...
2025-10-12 23:40:59,797 - INFO - Completed 3.txt: 131,130/131,701 valid
(99.6%)
2025-10-12 23:40:59,797 - INFO - [!] Crawl 0301 summary:
2025-10-12 23:40:59,797 - INFO - Total lines: 155,513
2025-10-12 23:40:59,797 - INFO - Valid records: 154,792
2025-10-12 23:40:59,797 - INFO - Malformed removed: 721
2025-10-12 23:40:59,797 - INFO - Duplicates removed: 0
2025-10-12 23:40:59,798 - INFO - [SUCCESS] Completed crawl: 0301
2025-10-12 23:41:00,002 - INFO - [CLEANUP] Cleaned up extracted files
2025-10-12 23:41:00,005 - INFO - [COMPLETE] Crawl data processing completed
successfully!
2025-10-12 23:41:00,005 - INFO - Processing time: 1312.8 seconds
2025-10-12 23:41:00,005 - INFO - Records processed: 904,874
2025-10-12 23:41:00,005 - INFO - Valid records: 898,254
2025-10-12 23:41:00,006 - INFO - Crawls processed: 2
2025-10-12 23:41:00,006 - INFO - Run validate_crawl_processing.py for detailed
validation and reporting

```

**Ingestion Speed:** ~280 videos/sec

#### Processed Data:

- 898,254 unique videos
- 2 crawls
- 19.1 million documents
- 3.16GB total storage
- ~53 minutes processing time

#### Query Performance:

- Simple lookups: <1ms
- Category aggregations: 50–200ms
- Cross-collection joins: 200–1000ms

#### Scalability:

- Linear for storage and processing
- 9M videos → ~9 hrs, ~31GB
- 90M videos → ~24 hrs, ~315GB (requires distributed processing)
- Logarithmic for query speed due to indexing

#### Efficiency Features:

- Incremental crawl processing (no reprocessing)
- 4-collection schema
- 1000-record batch upserts

- MongoDB sharding support
4. Source Code:
- Provide source code of your data preparation, data reduction, data transformation, and data ingestion steps in a ZIP file.

[Provided in the zip of the assignment submission.](#)