# Dremel Paper Review

**Name: Sharayu Rasal (N10802566)**

## Problem Statement Tackled

Google needed a way to run **interactive**, **ad-hoc queries** over **web-scale datasets** (trillions of rows, petabytes of data), much of it stored as **nested**, **non-relational records**. Traditional databases required slow loading phases and struggled with nested data at this scale.

**MapReduce** was excellent for batch processing but too slow for quick exploration and debugging. Dremel was built to **query data in place** (on GFS/Bigtable), support **nested structures directly**, and return **answers in seconds while scaling across thousands of machines**.

## Key Design Principles

- **Interactive speed**: Second-level responses for ad-hoc analysis, monitoring, and debugging.
- **In-situ access**: Operates directly on datasets in distributed filesystems (no heavy ETL or load step).
- **Columnar storage**: Stores data by column to scan only needed fields and compress better.
- **First-class nested data**: Preserves hierarchical structure without flattening (repetition/definition levels).
- **Massive parallelism**: Multi-level serving tree distributes work across thousands of CPUs.
- **Fault-tolerant by design**: Handles slow machines and failures i.e stragglers to keep latency low.

## Architecture & Components

- **Serving tree**:
  - **Root server** receives a query, reads table metadata, and rewrites the query into sub-queries.
  - **Intermediate servers** aggregate partial results flowing up the tree.
  - **Leaf servers** read tablets from storage and scan the required columns in parallel.
- **Query dispatcher** schedules work, balances load, detects stragglers via per-tablet timing histograms, and reschedules as needed.
  Supports returning results after a high percentage of tablets have completed when exactness trade-offs are acceptable.
- **Storage layout**: The tables are horizontally partitioned into tablets (typically replicated 3×).
  The leaf servers prefetch column blocks asynchronously to maintain high read hit-rates.

## Data Management (Nested Columnar Storage)

- **Column stripes for nested fields**: Values of a path like **Name.Language.Code** are stored contiguously, enabling selective scans.
- **Repetition levels**: Indicate where along a field's path a value repeats (e.g., Name or Language lists).
- **Definition levels**: Capture how many optional or repeated ancestors are present (including NULLs) to preserve structure.
- **Lossless structure**: Repetition + Definition levels allow precise reconstruction of nested records when needed
  Levels are bit-packed and NULLs are implicit.

## Query Processing

- **SQL-like language**: Path-based field references over nested data
  Also supports selection, projection, within-record aggregation, joins, top-k, and UDFs.

- **Execution on columns**: Select–project–aggregate advances column readers in lockstep and bypasses record assembly for speed.
- **Serving tree rewrite**: The root decomposes global aggregations into leaf partials, intermediates merge (e.g., sum of counts) and bubble results up.
- **Approximate operations**: One-pass algorithms for count-distinct and top-k when small errors are acceptable.

## Fault Tolerance & Scalability

- **Replication & failover**: tablets usually 3× replicated
  The leaves transparently switch replicas on read failures.
- **Straggler mitigation**: The dispatcher monitors tablet runtimes and redispatches outliers which allows <100% tablet completion can reduce tail latency when permitted.
- **Elastic parallelism**: near-linear speedups demonstrated from ~1,000 to 4,000 nodes
  User-perceived latency drops while total CPU time stays stable.

## Representative Use Cases at Google

- Analysis of crawled web documents.
- Android Market (Play) install tracking.
- Product crash reporting and spam analysis.
- Debugging Google Maps tiles.
- OCR results from Google Books.
- Datacenter job and disk I/O monitoring.
- Symbols/dependency analysis over Google's codebase.

## Related Work

- **MapReduce**: Great for large batch pipelines and fault tolerance but not tuned for low-latency ad-hoc queries.
- **Parallel DBMSs**: Strong at SQL and aggregation but historically row-oriented, flat schemas, and costly loading and is not commonly shown at thousands-node scale.
- **Dremel's niche**: It merges column stores (I/O/compression benefits)
   A nested data model (without flattening), and search-style serving trees (parallel fan-out/fan-in).
  Complements MapReduce (often querying MR outputs) and influenced BigQuery.

## Conclusion

Dremel makes **interactive analytics** feasible on **trillion-row**, **nested datasets** by combining
1. A nested columnar format (with repetition/definition levels)
2. SQL-style execution over columns (skipping row assembly for core operators)
3. A multi-level serving tree with robust dispatch and straggler control.

In practice, this delivers **order-of-magnitude speedups**
Map Reduce on rows → hours,
Map Reduce  on columns → minutes
Dremel → seconds.
Operating in situ on shared storage, it enabled rapid exploration, monitoring, and debugging at Google's scale, and set the blueprint for modern cloud analytics engines that must handle both scale and complex structure.