# INITIAL DRAFT

## Problem Statement Tackled

**Core problem:**

- Engineers need **interactive** (seconds-level) answers when exploring very large datasets (billions–trillions of records).

- Hardware & physics impose hard bottlenecks: disk I/O bandwidth and CPU cycles are finite; reading whole records unnecessarily wastes I/O and CPU.

- Data at web scale is often **nested** (lists, optional fields), not flat tables; flattening or fully normalizing nested data at petabyte scale is expensive or impossible in practice.

- Clusters are shared: tasks get preempted, nodes slow down (stragglers) or fail — any practical system must tolerate this.

So the problem Dremel solves: **return interactive, ad-hoc SQL-like query results over very large, read-only, nested datasets stored in place**, while (1) minimizing data read, (2) preserving nested structure without expensive reshaping, and (3) scaling to thousands of machines with straggler tolerance.

# Key Design Principles (derived from first principles)

I'll state each principle, then *why* it follows from first principles.

1. **Minimize I/O by reading only what's needed (columnar access).**

   o First principles: reading bytes from disk/network is expensive vs. CPU on modern clusters. If queries touch only a few fields, avoid reading unrelated fields.

   o Dremel stores each field's values contiguously (column stripes) so scanners can read only those bytes. This reduces disk bandwidth and enables much better compression (same-type data compresses better).

2. **Preserve nested structure without full materialization.**

   o First principles: nested data contains both values *and* structure (which element belongs to which list/record). If you drop structure you lose semantics. But storing whole records for every query is wasteful.

   o Dremel separates *structure metadata* (encoded as small integers: definition & repetition levels) from actual values so you can reconstruct arbitrary projections losslessly.

3. **Scale aggregation via a multi-level aggregation tree (serve like search engines).**

   o First principles: aggregating results from thousands of leaves to a single root is a communication/compute bottleneck at the root. Aggregation can be performed hierarchically to spread work and avoid a hot root.

   o Dremel uses a serving tree (root → intermediate → leaf) that rewrites queries into local partial aggregates which get combined up the tree. This reduces network fan-in and latency for many GROUP-BY/top-k queries.

4. **Work in-situ and interoperate with MapReduce-style pipelines.**

   o First principles: loading petabytes into a DBMS is costly and slows iteration. If the system can query files in their existing storage layout, analysts iterate faster.

   o Dremel reads data in-place (e.g., GFS/Bigtable) and complements — rather than replaces — MapReduce.

5. <mark>**Be robust to stragglers and multi-tenant interference.**</mark>

- First principles: on shared clusters, some tasks will be slow; waiting for every replica or slot can ruin latency.

- Dremel's dispatcher schedules tablets, uses replica failover, and lets you return results after a high-percentage of tablets complete (e.g., 98%) to trade slight accuracy for much lower latency.

6. **Avoid costly record (de)serialization where possible.**

- First principles: reconstructing typed objects from bytes (parsing) costs CPU; if queries can compute aggregates directly from compressed column values, you skip assembly cost. Dremel's internal iterators and evaluation algorithm operate on columns and levels, bypassing full record assembly for many queries.

# Architecture & Components (what's actually built and why it maps to the principles)

I'll list components, what they do, and the conceptual reason.

1. **Storage layer (GFS / distributed filesystem / Bigtable)**

   o Holds files/tablets in place (replicated). In-situ access removes load-time latency. Dremel reads compressed column stripes from these stores.

2. **Columnar nested storage: fields → stripes + levels**

   o **Column stripes**: contiguous runs of values for a field (A.B.C). Good for selective read and compression.

   o **Repetition levels (r)**: small integers attached to every stored value indicating *how deep* a repetition occurred (which repeated ancestor repeated last). This disambiguates which list element the value belongs to.

   o **Definition levels (d)**: small integers telling how many optional/repeated fields along the path *are actually defined*; used to represent NULLs and missing nested elements without storing explicit NULL markers.

   o *Why:* Together these two small integers encode record structure compactly and let you reconstruct arbitrary projections exactly. See the paper's sample records (r1/r2) and their columns for a concrete view.

*Mini example (from the paper, paraphrased):* field Name.Language.Code produces values ['en-us','en','en-gb'] across records; the repetition levels (0,2,1) and definition levels tell you which Name and Language instance each code belongs to — so you can reassemble nested records correctly.

3. **FieldWriter / DissectRecord (splitting algorithm)**

   o When writing/storing, a record is traversed and split into columns by a DissectRecord algorithm: child writers inherit parent levels and only update when they have data — this is a lazy, memory-efficient way to produce stripes with levels. (Appendix A in the paper gives pseudocode.)

4. **Finite-State Machine (FSM) for record assembly**

   o To reconstruct records from selected columns, Dremel builds an FSM where states correspond to field readers and transitions are labeled by repetition levels; traversing the FSM produces records with correct nested

structure. This enables *fast* and correct reassembly when needed. (Appendix B/C describe FSM construction and assembly.)

5. **Serving tree / query execution tree (root → intermediate → leaves)**

   o The root receives SQL-like queries and rewrites them into subqueries for intermediate nodes; leaves scan tablets in parallel and produce partial aggregates; intermediate nodes combine partials on the way up. This hierarchical aggregation reduces network and compute bottlenecks for GROUP-BY/top-k queries.

6. **Query dispatcher and slot model**

   o The dispatcher maps tablets to execution slots (threads), monitors per-tablet processing time histograms and reschedules stragglers. It supports parameters to return results after X% of tablets are scanned to allow trading completeness for latency.

7. **Internal execution within a server (iterators, select-project-aggregate evaluation)**

   o Dremel compiles optimized type-specific code for scalar functions. For many queries it uses a lockstep iterator algorithm (see Appendix D) that advances only the readers whose repetition level meets a fetchLevel and emits only expressions at appropriate nesting levels — this bypasses full record assembly and keeps execution one-pass and streaming.

8. **Prefetch + replication + read-ahead cache**

   o Leaves prefetch column blocks; a read-ahead cache yields ~95% hit rate in their setup. Tablets are usually 3× replicated; if a replica is slow/unreachable, leaves switch replicas. These techniques reduce read latency and improve robustness.

# How the important algorithms *work* (concise, stepwise)

(These are central to understanding the engineering.)

1. **How splitting a nested record into column stripes (conceptual steps):**

   o Walk the record top-down.

   o For each atomic leaf value encountered, compute:

      ▪ its *repetition level* (how many repeated ancestors have repeated since last occurrence), and

      ▪ its *definition level* (how many optional/repeated ancestors are actually present).

   o Append the value and its small r,d to that field's column block. If a repeated/optional field is absent, emit an implicit NULL via appropriate d value (no explicit null stored). This yields compact columns that can be scanned independently. (See DissectRecord pseudocode in Appendix A.)

2. **How reassembly via FSM works (conceptual steps):**

   o For the chosen subset of fields, construct an FSM describing which field reader to read next when a particular repetition level l is returned. The FSM encodes the nesting boundaries. (Appendix C gives the construction algorithm.)

   o Walk the readers following FSM transitions. When a value is non-NULL, move the output record's nested cursor to the appropriate level (start/end nested records as needed) and append the value. This reconstructs the exact nested structure for the selected columns.

3. **Select–Project–Aggregate evaluation without assembly:**

   o Keep a fetchLevel and selectLevel.

   o Advance only readers whose next repetition level ≥ fetchLevel. Compute next fetchLevel as the max next level. If WHERE holds, emit expressions whose repetition level ≥ selectLevel. This lockstep advancement means the engine can compute aggregates and scalar expressions directly on columns, avoiding object materialization. (Appendix D.)

# Related Work (where Dremel sits in the ecosystem)

- **MapReduce / Hadoop / Sawzall**: Designed for batch jobs — flexible, fault-tolerant but high latency. Dremel complements MR: MR is good for heavy transformations and long-running pipelines; Dremel is for interactive exploration of MR outputs.

- **Column stores (Vertica, MonetDB, etc.)**: Prior work shows columnar layout is efficient for analytical queries on *flat* relational data. Dremel extends columnar ideas to *nested* records using r,d metadata to encode structure.

- **XML compression/columnar techniques (e.g., XMill)** and **nested-relational algebra**: related in spirit (separate structure from content) but Dremel focuses on selective retrieval and interactive query performance at web scale.

- **Parallel DBMSs and hybrid systems (HadoopDB, etc.)**: Dremel emphasizes serving-tree aggregation and lockstep column processing to scale to thousands of nodes — a design point distinct from many parallel DBMS approaches at the time.

# Conclusion / Summary (key takeaways, tradeoffs, and practical guidance)

**What Dremel achieves**

- Interactive querying (seconds) over petabyte / trillion-row nested datasets by combining:

    o nested columnar storage (values + repetition & definition levels),

    o lockstep, column-oriented execution that often bypasses record assembly,

    o hierarchical serving-tree aggregation and robust dispatcher logic for stragglers and multi-tenant clusters.

**Quantified payoffs (paper's empirical findings)**

- Orders-of-magnitude improvements: e.g., in experiments switching MR from record format → columnar yielded ~10× speedup; using Dremel on top of columnar produced another ~10× improvement (paper's Figure 10). Dremel scanned terabytes on thousands of nodes in seconds for many queries.

**Tradeoffs / limitations**

- Record assembly and parsing remain expensive when you need full records — columnar wins for narrow projections; row (record-wise) storage can win when you need most fields. Dremel's crossover point varies by dataset.

- Stragglers and replication level matter: lower replication increases risk of slow queries; the system mitigates this by rescheduling and returning partial-but-accurate results when acceptable.

**First-principles, actionable guidance**

- If your workload is read-heavy, analytical, and often touches small subsets of fields (narrow queries), design for **columnar access + compact metadata** for structure.

- For nested data, you must encode structure cheaply (small integers per value) so columns remain compact and reconstructable.

- For aggregation at very large scale, hierarchical aggregation (multi-level tree) is the natural way to reduce fan-in bottlenecks.

- Always design for stragglers: rescheduling and replica choice are core controls for latency.

If you want, next I can:

- produce a diagram that visualizes the column stripes + r,d levels for the paper's sample records (hand-drawn ASCII or slide), or

- walk line-by-line through the DissectRecord pseudocode and show how it produces the r,d pairs for r1/r2, or

- generate a short notebook that *simulates* the repetition/definition encoding and demonstrates record assembly (toy Python).

Which of those would be most useful to you?

# DIVING DEEP IN "DATA MODEL" SECTION

Great — this is exactly the kind of close read that will make the rest of the paper easy. I'll do two things in this reply:

1. **Point out exactly which parts of my previous answer contain material that comes from the paper's *Data Model* section (so you can map back to what I already wrote).**

2. **Explain the Data Model itself in depth using first-principles reasoning**, with small, concrete examples you can work through by hand (so you'll be prepared for an exam on system architecture).

---

**1) Where the Data Model shows up in my previous answer**

The Data Model section of the paper (the part you pasted) is about the schema (Protocol Buffers), nested/typed records, multiplicities (required/optional/repeated), the tree structure and "on-the-wire" (row/sequence) representation — and the consequences (need to assemble records fast if you store values column-wise). Those ideas appear in these parts of my earlier answer:

- **Problem Statement Tackled**

  o "Data at web scale is often **nested** ... flattening or fully normalizing nested data ..." — this directly echoes the data-model motivation (nested, typed records).

- **Key Design Principles**

  o **"Preserve nested structure without full materialization."** — this point is directly driven by the schema properties (required/optional/repeated) and the need to keep structure when storing or querying.

- **Architecture & Components**

  o **"Columnar nested storage: fields → stripes + levels"** — that whole sub-section is the engineering realization of the Data Model: how to store Protocol-Buffer style nested fields as columns, and how to encode structure using small per-value metadata (what the paper later calls repetition/definition levels).

- o **"FieldWriter / DissectRecord (splitting algorithm)"** — the data model dictates how you must walk records to produce column values and the structural metadata.

- o **"Finite-State Machine (FSM) for record assembly"** — assembling the Protocol-Buffer style (row/sequence) view from columns is exactly what the FSM solves.

- **How the important algorithms work**

  - o **"How splitting a nested record into column stripes (conceptual steps)"** and **"How reassembly via FSM works (conceptual steps)"** — both explain, in algorithmic form, how the Data Model (the schema and wire layout) is converted to column storage and then reconstructed.

So: *everything in the paper that talks about "preserving nested structure", "splitting records into fields", "storing columns with small structure metadata", and "fast reassembly" is directly the engineering follow-on to the Data Model section you pasted.*

---

**2) Deep dive into the Data Model (first-principles + worked example)**

I'll start at the most basic questions and build up, so nothing is left hand-wavy.

**A. First principles: what problem does a data model solve?**

At the lowest level a system needs to represent information as bits on disk or on the wire. A *data model* is a formal vocabulary + structure rules that let programs agree on:

- what *kinds* of values can exist (integers, strings, floats), and

- how values are *composed* into bigger structures (records, lists), and

- what is required vs optional vs repeated, and

- how to serialize the result so different programs (C++, Java, Python, MR jobs) can read the same thing.

A good data model provides:

- **Type safety** (you know a field is an int),

- **Interoperability** (codegen in many languages), and

- **Compact on-the-wire representation** (efficient binary serialization).

Protocol Buffers (what the paper uses) is exactly that: *strongly-typed nested records* with required, optional, repeated multiplicities and a compact binary wire format. The Dremel paper's Data Model section just summarizes this.

## B. Concrete building blocks (from first principles)

1. **Atomic types** — the leaves of the structure: ints, floats, strings, booleans.
   First principle: different primitive types have different encoding costs and are the units of computation and compression.

2. **Record (or message / group) types** — named collections of fields (like structs).
   First principle: complex data = composition of primitives; compose them into named fields so you can address parts.

3. **Field multiplicity (semantics that matter!)**

   o **required**: the field **must** appear exactly once. If it's missing, the record is invalid. (Good for invariants.)

   o **optional**: the field may be absent (NULL / missing). You must be able to represent both "there and set" and "absent."

   o **repeated**: zero or more occurrences — i.e., a **list**. Order matters; repeated fields are sequences.

First principle: these multiplicities are *semantic* — they change how you interpret presence/absence and how you store and reconstruct values.

4. **Tree / path notation**

   o Nested records form a *tree*. Each leaf is addressed by a *path* (dotted notation) like Name.Language.Code.
   First principle: any nested field can be uniquely identified by its path; that lets you reason about projections (selecting a subset of paths).

5. **On-the-wire serialization**

   o Protocol Buffers (as described in the paper) serialize records **sequentially as they occur** in the record (row-wise). That is, when you write a record, you walk the record and append each field's encoded bytes in order. This yields a compact, language-neutral binary representation.

   o First principle: row-wise (sequential) serialization is ideal for streaming a record to another program; it is natural and efficient for producers and consumers that work record-by-record.

## C. Why that row/sequence representation matters for a columnar system

Dremel's goal is interactive analytics on very large nested datasets. Two conflicting facts:

- Protocol Buffers (row/sequence) is great for cross-language producers/consumers and MR jobs that expect records.

- Columnar storage (one file/stream per *path*) is great for analytics: when a query touches only two fields out of 100, reading only those columns saves large I/O.

**Consequence (first principle):** If you store data column-wise but other tools expect row/sequence (protobuf) records, you must be able to **reconstruct exact records quickly** from the columns. That reconstruction is what forces Dremel to design compact *structure metadata* (the repetition/definition level idea) and fast assembly algorithms (DissectRecord, FSM).

### D. How Dremel maps the tree model to columns (conceptual picture)

Think of the schema as a tree of fields. For each *leaf path* (every place where there is an atomic value) Dremel stores a **column** — a contiguous list of the values that occurred for that path across the dataset. But just storing values loses the structure (which values belong to which record, which values are missing, and which values are elements of repeated lists).

So Dremel stores, **alongside each value**, two tiny integers that encode structure:

- **Definition level (d)** — encodes how many fields *along the path* are actually defined (not missing). This answers questions like "is the optional ancestor present?" or "is this value a NULL?" Conceptually it tells you how deep the value goes in terms of optional presence.

- **Repetition level (r)** — encodes list boundaries: it tells you whether this value continues the same repeated container(s) as the previous value, or starts a new list element (or even a new top-level record). Conceptually it tells you "which repeated ancestor(s) did we continue vs. start over."

These two small integers per value *let you reconstruct* the hierarchical record shape precisely from the purely value-centric column.

First principle recap: we traded a few small ints per value for the ability to store values densely and still rebuild the original nested records — that trade gives huge savings in I/O and compression while retaining fidelity.

### E. Walkthrough example (small, concrete, by hand)

Schema (plain, readable form):

Document {

 required int DocId;

 optional group Links {

```
    repeated int Forward;

    repeated int Backward;

  }

  repeated group Name {

   repeated group Language {

    required string Code;

    optional string Country;

   }

  }

}
```

Notes:

- Links is optional (may be missing).

- Forward and Backward inside Links are repeated lists.

- Name is repeated; Language inside Name is repeated. Each Language has required Code and optional Country.

Two sample records (row order):

**r1**

DocId: 1

Links:

 Forward: 10, 11

Name:

 - Language:

   - Code: "en-US", Country: "US"

   - Code: "en"        (no Country)

 - Language:

   - Code: "fr-FR", Country: "FR"

**r2**

DocId: 2

(no Links)

Name:

 - Language:

   - Code: "es-ES", Country: "ES"

Now focus on one leaf path: Name.Language.Code. If we walk the records in row order and emit Code values as we encounter them, the **sequence of values** is:

["en-US", "en", "fr-FR", "es-ES"] — values for r1 followed by r2.

If we store *only* these values in a column, we cannot tell:

- which codes belong to which Name element,

- which Language elements belong to the same Name,

- or where r1 ends and r2 begins.

**So** Dremel stores, with each value, small metadata indicating "how defined" and "how repeated". I'll explain conceptually how to read it:

- For value "en-US" (first in r1): this is the *first* Language inside the *first* Name of *r1*. So we would mark it as starting a new record and starting new repeated containers → a low r (meaning no repeated ancestor is being continued from previous value) and a d that marks Code is present (and Country present for this value).

- For value "en" (second in r1): this is another element of Language **within the same** Name as the previous value. So the repetition metadata must indicate "same Name, continued Language list" (i.e., it's another Language entry; r is set so assemblers know not to start a new Name node).

- For "fr-FR" (third r1): this belongs to the *second* Name record in r1, so r must indicate the boundary that tells the reassembler "finish the first Name, start the second Name, start its first Language".

- For "es-ES" (r2): this is in a new top-level record — so the metadata must indicate "new Document".

The reassembler (FSM) uses those markers to open/close nested objects and place each value at the correct path.

**Important conceptual point:** you can think of r and d as tiny breadcrumbs that tell you two things:

- d: "how many ancestors are actually present (not NULL) for this value?"

- r: "which repeated ancestor levels are we continuing vs starting anew compared to the previous value?"

(You do not need to memorize numeric formulas for r and d in this moment; understand their *meaning* — exam questions will often probe that you know what problem they solve and how they are used to re-make records.)

**F. Sketch of the algorithms the data model forces you to write**

1. **Writer / Dissection (producing columns from records)** — what it must do (conceptual steps):

   o Walk the record tree depth-first.

   o Each time you hit a leaf value, append it to that leaf's column *and* compute:

      ▪ the definition level — based on which optional ancestors are present for that value,

      ▪ the repetition level — based on whether repeated ancestors are being continued or are starting fresh.

   o If a repeated field is empty, you still encode that by emitting the appropriate definition information (so readers know the list is empty).

   o Output is: per path (column) → a list of (value, d, r) tuples stored in order encountered.

2. **Reader / FSM (reconstructing rows from columns)** — what it must do (conceptual steps):

   o For the selected set of columns, create an FSM that expresses nesting boundaries (the schema tree restricted to those columns).

   o Advance columns in lockstep, reading the next (value, d, r) from whichever column produces the next low-level event, and use r/d to decide:

      ▪ whether to close or open repeated lists or optional containers,

      ▪ whether to insert NULLs for missing optional fields,

      ▪ how to place the value in the current output record.

   o Emit completed records in the same order as original rows.

3. **Why the FSM is necessary**

- Because columns are independent streams of values; each value carries only local context (r,d). The FSM coordinates those local contexts into global structure events (start Name, end Name, start Language, etc.).

## G. Why the paper emphasizes "code generation" and the wire format

The Data Model section mentions that code generation produces bindings for C++/Java and that records are serialized sequentially. That is important because:

- Many tools (MR jobs, analytics scripts) expect a record/row view or expect protocol buffers objects. If Dremel kept data only as exotic column files, those tools would break.

- So Dremel must make it cheap to **construct** the Protocol-Buffer style record object from columns when needed. That's why a fast assembly strategy (FSM + compact levels) is vital: it preserves interoperability with the rest of Google's ecosystem (and with external MR programs).

## H. Exam-ready summary of what to memorize / understand deeply

- **Data model = Protocol Buffers style nested, strongly-typed records.** Know what required, optional, repeated mean semantically.

- **The logical schema is a tree; every leaf path can be treated as a column.** Dotted notation A.B.C addresses leaves.

- **Wire format is row/sequence oriented.** That matters because columnar storage needs to provide a way to reconstruct that row view.

- **Dremel's trick:** store each leaf column as a dense sequence of values *plus* tiny structure markers (definition + repetition levels) for each value. That lets you reconstruct exact rows without storing full records or expensive pointers.

- **DissectRecord (writer) and FSM (reader)** are the components that convert between row (protobuf) and column representation. On an exam, you should be able to explain their job and the *why* behind each step.

- **Big picture:** the Data Model causes two design constraints for the system: 1) preserve semantics of optional/repeated/required; 2) allow efficient projection (read only the needed columns). Solutions (r/d + FSM) are the direct answers to those constraints.

---

### Final notes & suggested next steps (I can do these right now if you want)

- I can **draw a diagram** of the sample schema + columns + annotated (value, d, r) streams so you can visually trace assembly.

- I can **simulate the encoder + decoder** in a tiny Python notebook and show the exact (value, d, r) numbers for the example above and run through the FSM to reconstruct r1 and r2 (this is an excellent way to internalize the semantics).

- I can **walk line-by-line through the DissectRecord pseudocode** from the paper and show exactly where the d and r values are computed (appendix A of the paper).

Which would you prefer next? Diagram, runnable toy simulation, or a step-through of the paper pseudocode?

# DIVING DEEP IN "NESTED COLUMNAR STORAGE"

Nice — let's do the same focused, first-principles treatment for **"Nested Columnar Storage"** (Section 4). First I'll point to which parts of my earlier answer already covered this material, then I'll explain every idea in the section in depth, with concrete reasoning and worked examples so you can nail a systems-architecture exam. (I'll reference the paper where appropriate.)

---

**Where this section appeared in my earlier answer**

The paper's *Nested Columnar Storage* section maps directly to these parts of my previous reply:

- **"Columnar nested storage: fields → stripes + levels"** in Architecture & Components — that paragraph summarized repetition/definition levels, columns, and reconstruction.

- **"FieldWriter / DissectRecord (splitting algorithm)"** — this is the same as the Splitting Records into Columns subsection (4.2).

- **"Finite-State Machine (FSM) for record assembly"** and the algorithmic sketches — these correspond to Record Assembly (4.3).

So what follows is a *deeper* walk-through of Section 4 — rationale, mechanics, and examples — from first principles.

---

**First principles recap (why nested columnar storage is needed)**

1. I/O is expensive at scale. Reading fewer bytes → faster queries.

2. Analytical queries usually touch a small set of fields (projection), so storing each field's values contiguously (columnar) saves I/O and compresses better.

3. But nested records are *not flat*: lists and optional fields make the mapping from values → records ambiguous if you only store values.

4. Therefore: we must store *compact structural metadata* alongside values so we can reconstruct the exact nested records losslessly while still reaping column I/O benefits.

Section 4 provides that metadata design (definition + repetition levels), describes how to compute/encode them efficiently, and explains how to reconstruct records.

**1) Repetition & Definition levels — what they are and *why* they solve the problem**

**What problem each solves (intuitively)**

- **Repetition levels (r)** answer: *When the next value comes, which repeated ancestor did we repeat at?*
  → They disambiguate whether successive values belong to the same list element, a new list element in the same parent, or an entirely new top-level record.

- **Definition levels (d)** answer: *For this value (or NULL), how many optional/repeated ancestors along the path are actually present?*
  → They let us represent NULLs / missing values compactly and tell us which parent containers exist.

Think of each emitted leaf value as carrying two tiny breadcrumbs: one that says "how deep the present structure goes" (definition), and one that says "which repeated boundary was crossed compared to the previous value" (repetition). With these two pieces we can rebuild the tree.

**Formal constraints (first-principles)**

- For a leaf path p, let R be the number of repeated ancestors on that path. Then repetition level r ∈ {0, 1, ..., R} where:

  - r = 0 denotes start of a new top-level record (the record boundary).

  - larger r denotes repetition at deeper repeated ancestors (higher numbers = deeper repeated fields repeated since previous value).

- Definition level d ∈ {0, 1, ..., Dmax} where Dmax equals the number of fields on the path that *could be undefined* (optional or repeated) — d counts how many of those are actually present at this value.

  - If a leaf value is missing (NULL), we still emit a metadata entry (a NULL marker) whose d indicates how far the structure was defined.

Why integers (levels) instead of is-null bits?

- Using integer levels captures *which ancestors are present*, not just "defined or not." That extra information is necessary to reconstruct nested parents (e.g., whether Name exists even when Country is NULL). It also avoids storing explicit NULL tokens and allows compact bit-packing of small integers.

**Concrete small example (re-using the paper's running schema)**

Schema path: Name.Language.Code.

- Name is repeated, Language is repeated, Code is required. So R = 2 (two repeated ancestors). Repetition levels range 0..2.

- For r1 (paper example), Code occurrences and repetition levels: en-us → r=0, en → r=2, en-gb → r=1. This sequence tells the assembler which repeated container (Name or Language) changed between successive values.

Why a NULL is inserted between en and en-gb in paper's example?

- Because the second Name had no Code values at all (i.e., that Name had no Language with a Code) — to represent that absence in the column stream we insert a NULL entry with appropriate d. That lets the assembler know there was an intermediate Name with no Code children.

---

## 2) Encoding — how levels and values are stored

### Block structure

- Each column is stored as a sequence of **blocks**. Each block contains:
  - compressed field values (only actual non-NULL values are stored),
  - the packed repetition and definition levels aligned with those values (and including entries for NULLs as needed).

- NULLs are **not** stored as values — they are implied by definition levels smaller than the maximum. So column value buffer contains only real values; logical NULL slots are inferred when reading definition levels.

### Compactness & bit-packing

- Maximum d or r for a column is small (path depth rarely > 8 in practice). So levels fit in a few bits each. The implementation uses just enough bits to encode the maximum observed level, then packs levels densely.

- Repetition levels are omitted when unnecessary (e.g., def level 0 implies r=0, so r can be skipped). Similarly, if a field is always defined (no optional or repeated ancestors that could be undefined), you can omit definition levels.

### Why this is smart (first principles)

- We minimize bytes per value: the value itself is compressed; structure metadata is tiny and bit-packed; NULLs do not consume value storage. This preserves the columnar I/O benefit while retaining exact reconstructability.

---

## 3) Splitting Records into Columns (4.2) — algorithms & implementation pattern

**Problem**

Given a nested record, produce per-path columns containing (value, d, r) streams in a single pass and cheaply handle missing fields for sparse schemas.

**High-level idea (first-principles)**

- Walk the record tree once (depth-first). When you encounter an actual atomic value, compute its d and r and append it to the value stream for that leaf.

- But many fields are missing in sparse datasets; we must avoid walking/store work for thousands of unused fields for each record.

**Field writer tree (practical trick)**

- Create a tree of **field-writer** objects mirroring the schema. Each leaf writer is responsible for a column. Writers *inherit* levels from their parent; they only update (synchronize) parent state when they actually produce data.

- Child writers don't force their parents to produce entries unless a child needs to write — this lazy propagation avoids touching thousands of empty columns per record.

Concretely:

1. For each incoming record, start at the root writer with current levels set to "start of record".

2. When you reach a field with a concrete value:
   - compute local d and r using parent's current levels and the schema path,
   - append the value to that field's column block along with the level bits.
   - mark children accordingly if they exist.

3. If a repeated field is empty, the algorithm must still emit appropriate definition-level entries (logical NULLs) for leaf columns that are missing because of this omitted subtree. The writer tree knows how to emit the needed NULL markers cheaply (no value bytes — only small packed level entries).

Why child writers inherit parent state?

- Because most of the time you traverse a nested value that shares many ancestor definitions with the previous leaf; inheriting avoids recomputing higher-level status and lets you set levels incrementally.

**Why this is efficient (first principles)**

- Work per record is proportional to number of present fields, not the full schema size. For sparse, wide schemas (many possible fields, few present), this saves huge CPU & memory. Also, because writers append into contiguous blocks, compression and I/O are efficient.

---

## 4) Record Assembly (FSM) — reconstructing rows from columns (4.3)

### Goal

Given a subset of columns, reconstruct records *as if* the original records existed but stripped of non-selected fields — preserving nesting and order.

### Key idea

- Build a **finite state machine (FSM)** whose states are associated with selected field readers. Transitions are labeled with repetition levels.

- Each field reader yields (value, d, r) in order. The FSM consumes these and emits nested structure events (start/stop repeated/optional containers and leaf assignments) to create output records.

### How transitions are chosen (concise)

When the current reader f reports its *next* repetition level l, we:

1. Starting from f in the schema tree, find the *ancestor that repeats at level l*.

2. Select the **first leaf field** inside that ancestor — call it n.

3. Create a transition (f, l) → n in the FSM: if the next repetition level is l, we should next read values from n.

This rule ensures that when you see an r indicating a repetition at some ancestor boundary, you jump to the correct next field in record order.

### Worked FSM example (paper's example paraphrased)

- Start state: DocId.

- After reading DocId, FSM transitions to Links.Backward. The backward values are drained (possibly zero), then FSM goes to Links.Forward. After those, it goes to the Name.* subtree and so on.

- If reading Name.Language.Country and the next r = 1, the algorithm finds that the ancestor with repetition level 1 is Name. The *first leaf* inside Name is (paper example) Name.Url. So the FSM transitions to Name.Url.

This yields correct sequential reconstruction: the FSM will traverse once per record and append values (or NULLs, inferred from definition levels) in the correct nested places.

**Why FSM rather than ad-hoc merging**

- Columns are independent ordered streams. Coordinating them requires a systematic rule to decide next reader so you preserve global record order and nesting. FSM with r-labeled transitions encodes that rule compactly and runs in linear time over the number of emitted values/levels.

---

**5) Example: compute (value, d, r) stream for the sample records (walk-through)**

Using the same schema & records as earlier (Document with DocId, optional Links, repeated Name → repeated Language with Code required and Country optional). We'll compute for two leaf columns: Name.Language.Code and Name.Language.Country. (These are exactly the fields discussed in the paper's figures.)

**Step A — identify repeated and optional ancestors**

- For both Code and Country:
  - Name is repeated → contributes to repetition depth.
  - Language is repeated → contributes to repetition depth.
- For Code:
  - Code is **required** in Language, so Code cannot be NULL when a Language exists.
  - Thus definition-levels for Code tell us whether Language exists (and whether Name exists).
- For Country:
  - Country is **optional**, so Country may be absent even if Language exists. So its definition levels encode finer distinctions.

**Step B — write out values and label r,d (as paper does)**

From paper's narrative, for Name.Language.Code in r1:

- values: en-us, en, en-gb
- repetition levels: 0, 2, 1 (explained earlier).
- Between en and en-gb, there is an empty second Name (no Code values), so we insert a NULL entry with a suitable definition level to represent that missing Language in that Name.

For Name.Language.Country:

- Some occurrences exist (with country) and several are missing. The definition levels for missing entries differ depending on whether the parent Language exists but Country omitted (higher d) or the whole Language is missing (lower d).

- Concretely (paper): missing occurrence of Name.Language.Country in r2 has d = 1 (meaning only Name exists? or maybe that Name is defined but Language missing) while missing occurrences in r1 have d = 2 (inside Name.Language) and d = 1 (inside Name). The point: d encodes how many optional/repeated ancestors are present.

I'm intentionally focusing on *meaning* rather than forcing you to memorize numeric conventions — but for exam clarity:

- Compute R = number of repeated ancestors and Dmax = number of optional+repeated ancestors that can be undefined.

- For each emitted value or NULL, $r \in [0..R]$ and $d \in [0..Dmax]$. Higher d means deeper ancestors exist.

If you want, I can produce a **concrete table** of (value, d, r) for each emitted element (including inserted NULLs) using the exact r1/r2 records; say the FSM and the paper already do that in Figure 3 — I can reproduce a similarly explicit table for your study guides.

---

**6) Practical implementation notes & optimizations (what to expect in real systems)**

- **Omit storing levels when constant.** If a field is always defined, skip definition levels. If repetition levels implied by d=0, skip them. This reduces metadata overhead in common simple columns.

- **Blocks & compression:** pack levels densely per block; compress values separately; blocks are the unit for prefetch/read-ahead.

- **Sparse schema optimization:** field-writer tree is essential — avoid per-record work proportional to schema width.

- **Partial projections:** when reading a subset of fields, build a reduced FSM over only those paths — much cheaper than full reconstruction.

---

**Exam checklist: what to memorize vs. what to reason about**

Memorize:

- Definitions & roles of **repetition** and **definition** levels (what question each answers).

- The high-level pipeline: **DissectRecord → columns with (value + packed levels) → FSM assembly**.

- That NULLs are *not* stored as values; they are encoded via definition levels.

- That repetition levels range 0..R where R is the number of repeated ancestors; 0 marks a new top-level record.

Be prepared to **reason** (not memorize):

- Given a simple nested schema and two consecutive column values, explain how r indicates whether you stay in the same repeated container or start a new one.

- Given a missing leaf, explain how a d value disambiguates a missing leaf vs. missing parent container.

- Explain why storing small integers per value is cheaper than storing full pointers or duplicate structure.

---

**Next steps I can do immediately (pick any; I'll do it now)**

- Produce a **table** that lists every emitted entry for Name.Language.Code and Name.Language.Country for r1 and r2, with explicit (value, d, r) including the NULL rows (great for memorization).

- Draw a **diagram** showing the schema tree, column streams, and how the FSM traverses them for one record (visual aid).

- Generate a **tiny Python simulation** that takes the small sample records and emits column streams with levels, then uses a simple FSM reassembler to reconstruct the records (hands-on and unbeatable for exam prep).

Which of those would you like me to generate right now?

# DIVING DEEP IN "QUERY LANGUAGE & EXECUTION"

Great — let's break down **Query Language** and **Query Execution** exactly like before: first I'll point to where these topics appeared in my earlier summary, then I'll unpack every idea from first principles, give concrete examples and pseudo-algorithms you can use in an exam, and finish with an exam checklist and suggested next steps (diagrams / runnable sims I can produce right now).

No clarifying questions — I'll give you a best-effort, exam-ready deep dive.

---

**Where this material appeared in my earlier answer**

Parts of my previous reply that already cover the content from the *Query Language* and *Query Execution* sections:

- **Architecture & Components**

  - The **serving tree / query execution tree (root → intermediate → leaves)** description maps to the Tree Architecture and query rewrite/aggregation idea in Query Execution.

- **How the important algorithms work**

  - **Select–Project–Aggregate evaluation without assembly** — this explains lockstep iterators, fetch/select levels, and bypassing full record assembly during execution.

- **Key Design Principles**

  - Points about minimizing I/O, doing in-place reads, hierarchical aggregation, and straggler tolerance are all core to Query Execution.

So what follows drills into the semantics of the language and the actual runtime mechanics the paper describes.

---

**QUERY LANGUAGE — first principles + deep explanation**

**1) High-level intent (first principles)**

- You want a language that feels like SQL because analysts know SQL and it expresses projection, selection, aggregation naturally.

- But the underlying storage is *columnar* and *nested*. So the language must:

    1. Express operations over **nested** tables using path expressions (e.g., Name.Language.Code).

    2. Have semantics that are implementable **efficiently** on columnar nested storage (so we can exploit column pruning, skip reading irrelevant columns, and compute aggregates in one pass).

- Dremel chooses a SQL-like language with path expressions and semantics adapted to nested inputs — it keeps SQL familiarity but defines how projection/aggregation map onto nested structures.

**2) Key semantic ideas (and why they matter)**

A. **Path expressions**

- Use dotted paths (A.B.C) to address fields inside nested records.

- First principle: a nested schema is a tree; dotted paths uniquely identify leaves. This allows the query planner to pick the exact columns to read.

B. **Selection (WHERE) on nested records**

- Think of a nested record as a **labeled tree**. A selection predicate is interpreted as a pruning operation on that tree: branches that don't satisfy predicates get removed.

- Example: WHERE Name.Url LIKE 'http%' removes Name subtrees whose Url doesn't match; other Names in the same Document that do match are retained.

- Why this is important: selection can be applied *per-subtree*, not just per-row, enabling within-record filtering.

C. **Projection placement & repetition semantics**

- **Rule:** Each scalar expression in SELECT emits a value at the nesting level of *the most-repeated input field used* in that expression.

    o Intuition: If you compute concat(Name.Language.Code, '-'), and Name.Language.Code is repeated twice deeper than DocId, the resulting concatenation is naturally an element inside the same repeated container — so place it at that repeated level.

- Why: This rule preserves consistent nesting in results without requiring the user to write explicit record constructors. The result's schema is derived from input multiplicities.

D. **Within-record aggregation**

- Aggregates can be done **within a subrecord** (e.g., count number of Language.Code per Name).

- Example: SELECT Name.Url, COUNT(Name.Language.Code) — the COUNT value is computed for each Name subtree.

- First principles: Aggregation groups are defined by the nearest repeated ancestor you want the result attached to. Dremel supports GROUP BY and within-record aggregations (and these are implemented efficiently in the serving tree as partial aggregates).

E. **Supported features**

- SQL-like: projections, selections, GROUP BY, nested subqueries, joins, top-k, UDFs, approximate algorithms (top-k, count-distinct) for one-pass execution.

**3) Concrete example walk-through (simple)**

Input table T contains documents r1, r2 (paper's example). Example query:

SELECT

 Name.Url,

 uint64(COUNT(Name.Language.Code)) AS num_codes,

 CONCAT(Name.Language.Code, '-') AS code_tag

FROM T

WHERE Name.Url LIKE 'http%'

Explain:

- WHERE prunes Name subtrees per-document: only Names with Url starting with http survive.

- COUNT(Name.Language.Code) is computed *per Name* (within-record aggregation), returning one integer per Name.

- CONCAT(…Code…) emits a value at the same level as Name.Language.Code (a deep repeated level), so for each language code there will be an associated code_tag.

- Result: a nested output where each retained Name contains num_codes (one integer) and a sequence of code_tag values.

**Exam point:** Be able to explain how selection prunes tree branches, how aggregate attaches to the right repeated ancestor, and how expression placement is decided by the deepest repeated input field.

**QUERY EXECUTION — first principles + deep explanation**

**1) High-level architecture: the multi-level serving tree**

- **Problem to solve:** Queries over huge tables must scan many data shards (tablets). Sending all raw data to one node kills network and compute; we must parallelize and aggregate partially to avoid bottlenecks.

- **Solution (first principles):** Build a **serving tree** (root → intermediate → leaf nodes) so that:

  - Leaves scan local tablets in parallel and produce *partial results*,

  - Intermediate nodes combine partials (e.g., SUMs of counts per key),

  - Root finalizes the aggregation and returns a small-to-medium-sized result.

This is the classical *divide-and-conquer* rewriting of aggregation over partitions:

- For SELECT A, COUNT(B) FROM T GROUP BY A:

  - Leaf i computes SELECT A, COUNT(B) AS c FROM T_i GROUP BY A.

  - Intermediate nodes do SELECT A, SUM(c) FROM (union of child results) GROUP BY A.

  - Root final SUM finalizes into global per-A counts.

**Why this is good:** It reduces network fan-in and distributes CPU (partial aggregation reduces data volume moving up the tree).

**2) Query dispatcher & scheduling (slots, histograms, straggler mitigation)**

- **Slots:** The system has a fixed number of execution slots (threads across all leaf nodes). Example: 3,000 leaves × 8 threads = 24,000 slots.

  - Think of slots as units of concurrency for scanning tablets.

- **Tablet assignment:** There are often more tablets than slots, so each slot processes multiple tablets (e.g., ~5 tablets/slot). The dispatcher maps tablets → slots.

- **Histogram & rescheduling:** During execution, dispatcher builds a histogram of tablet processing times. If a tablet is unusually slow (a straggler), the dispatcher reschedules that tablet (or assigns another replica) to another slot — repeated as necessary. This reduces tail latency.

- **Replica fallback & read-ahead:** Tablets are replicated (commonly 3x). If a leaf server can't reach a replica or a replica is slow, the leaf falls back to another replica. Blocks in stripes are prefetched asynchronously; read-ahead caches achieve ~95% hit rate — this minimizes I/O wait.

- **Minimum-percent parameter:** Dispatcher can be configured to require only X% of tablets to be scanned before returning results (e.g., 98% instead of 100%).

  - **Why:** For many aggregations, missing a few tablets (rare) changes answers negligibly; for interactive latency, returning early with approximate results is valuable. Systems can trade perfect completeness for faster responses.

**Exam-ready pseudo-code for dispatcher behavior (conceptual):**

for each incoming_query:

 tablets = list_tablets(query.table)

 assign tablets to slots (spread across leaves)

 start scanning tasks

 while unfinished:

  collect tablet_completion_times into histogram

  if any tablet_time >> median * threshold:

   reschedule tablet on another replica/slot

  if percent_complete >= min_percentage:

   if query_allows_early_return:

    gather partials and return aggregated result

 gather all partials -> combine up serving tree -> return final result

**3) Internal execution on a server — lockstep iterators & bypassing assembly**

- **Key idea (first principles):** For many analytics queries you don't need to reconstruct full records. If you can compute aggregates and evaluate expressions while scanning columns, you save the cost of materializing objects and re-parsing. Dremel implements **lockstep iterators** that advance selected column readers in coordination by using repetition & definition levels.

- **FetchLevel & SelectLevel mechanism (compact):**

  - fetchLevel determines which readers must advance to produce the next structural event.

- selectLevel determines the nesting granularity at which results and expressions should be emitted.

- The iterator algorithm reads only the necessary columns, advances readers whose repetition levels meet the fetchLevel, evaluates predicates, and emits aggregate contributions annotated with correct repetition/definition levels — all without building row objects.

- **One-pass aggregates:** For SUM/COUNT/MIN/MAX you can compute partial aggregates in streaming fashion. GROUP BY requires grouping keys; the serving tree typically groups by keys at leaf level (or uses local hash tables) and produces partial grouped results.

**Conceptual lockstep iterator pseudocode (simplified):**

initialize readers for required columns

while not all readers exhausted:

  determine next fetchLevel (max rep level among nexts)

  advance readers with rep >= fetchLevel

  evaluate WHERE predicate using current reader values & def levels

  if predicate true:

    compute expressions / update aggregates

  when group boundary reached:

    emit partial aggregate with repetition level info

- **Optimized scalar code generation:** Dremel compiles tight, type-specific code for scalar functions (e.g., string concat, arithmetic) to speed up execution — less dynamic dispatch, more inlined code.

**4) Serving tree partial aggregation rewriting (formalized)**

- Given a global aggregation AGG_GLOBAL = AGG over T, rewrite as hierarchical aggregation:

Let T = union of $T_i$ (tablets or partitions).
Leaf queries: $R_i$ = SELECT key, AGG_local(value) FROM $T_i$ GROUP BY key
Intermediate/root: SELECT key, AGG_combine($c_i$) FROM union($R_i$) GROUP BY key
Where AGG_combine is the combining operator (e.g., SUM combines counts, MIN picks min, AVG uses (sum, count) pair).

This is the standard *monoid* decomposition idea: any associative combiner lets you do partial aggregation safely.

**Exam tip:** Be able to show the rewrite step for COUNT, SUM, MIN, and how AVG requires two partial columns (sum, count).

## 5) Handling different query classes

- **One-pass aggregations & projections** — Dremel excels: single scan, local partials, hierarchical combine.

- **Top-k / count-distinct** — Dremel may use approximate streaming algorithms (e.g., HyperLogLog, lossy-counts) that fit one-pass and limited memory.

- **Large aggregations or joins** — might need more heavyweight engines (parallel DBMS techniques, external shuffles) or MapReduce-style pipelines. Dremel is optimized for interactive, small-to-medium result sets.

## 6) Fault tolerance, partial results, and tradeoffs

- **Rescheduling & replication** provide fault tolerance and straggler mitigation.

- **Early return on <100% scanned** offers latency improvement at cost of exactness — useful for exploratory analytics.

- **Memory vs I/O tradeoffs:** Partial aggregates may require local memory (hash tables). If memory inadequate, system must spill or use multi-pass algorithms.

---

## Concrete worked example of query execution (end-to-end)

Query:

SELECT Name.Url, COUNT(Name.Language.Code) AS c

FROM Documents

WHERE Name.Url LIKE 'http%'

GROUP BY Name.Url

Execution sketch:

1. Root: parses query, finds tablets spanning Documents.

2. Root rewrites to intermediate: each intermediate server gets subset of tablets and rewrites to leaves.

3. Leaves: for each tablet T_i, leaf executes:

   o Use column readers for Name.Url and Name.Language.Code.

   o Lockstep iterate, apply WHERE (prune Name subtrees where Url not like 'http%').

- o For each Name that passes, compute COUNT(Code) locally (hash by Name.Url if many Names per tablet).

  - o Emit pairs (Name.Url, local_count) up to parent.

4. Intermediate: sum partial counts per Name.Url from its children.

5. Root: sum across intermediates, produce final Name.Url, total_count list.

6. If min_percentage < 100% and enough tablets done, root may return intermediate result early.

---

**Exam checklist: facts & reasoning to memorize**

- Query language:

  - o Path expressions for nested fields.

  - o WHERE prunes branches; projection emits values at level of most-repeated input field used.

  - o Within-record aggregation attaches aggregates to the repeated ancestor that defines the grouping.

  - o Support for nested subqueries, top-k, joins, and approximate algorithms.

- Query execution:

  - o Multi-level serving tree: leaves → intermediate → root; partial aggregation at each level.

  - o Dispatcher: slots, histograms, rescheduling stragglers, replica fallback, min-percent scanned parameter.

  - o Lockstep iterators: read selected columns in coordination using repetition & definition levels; bypass assembly for many queries.

  - o Aggregation rewrite: local partials + combining function at higher levels (monoid property).

  - o Approximate algorithms (top-k, count-distinct) are used when single-pass/low-latency is required.

- Practical tradeoffs:

  - o Early results vs completeness; memory for local grouping; complexity for joins or very large results.

---

**Suggested next steps (pick one — I'll do it immediately)**

1. Produce a **diagram** of the serving tree showing how a GROUP BY query is rewritten and aggregated up the tree.

2. Generate a **detailed end-to-end example** (table of leaf outputs, intermediate aggregates, final result) for the small query above.

3. Produce **runnable pseudocode** (Python-like) that simulates:

   - a) dispatcher scheduling + rescheduling on stragglers, and

   - b) a simplified lockstep iterator that reads precomputed (value,d,r) streams and computes a per-Name COUNT.

Which would you like?