# On-Demand Schema Migrations for Multi-Model Data Integration: A Topos-Theoretic and Generative AI Perspective

Matthew Long

`Magneton Labs`

January 22, 2025

**Abstract**

Modern enterprises often store their data across heterogeneous systems: relational databases, NoSQL document stores, and graph databases. This diversity poses challenges for schema migrations, especially when queries or analytics must draw on data stored in incompatible models. In this paper, we detail a topos-theoretic perspective on schema migrations, demonstrate on-demand transformations that allow users and applications to query data regardless of its storage model, and highlight how generative AI agents can leverage these concepts to unify multi-model data seamlessly. We offer mathematical formalisms for data-model transformations, show the benefits of this approach, and present a Haskell-based sample implementation for a schema migration engine.

## 1  Introduction

As data proliferates across relational, document, and graph databases, organizations face persistent challenges in aligning or migrating schemas. Often, the need arises dynamically during query execution: a user or application attempts to fetch data that resides in a different representation, potentially with incompatible schemas. Traditionally, schema migrations are large, planned events. Here, we argue for a more fluid, on-demand approach based on category and topos theory.

Moreover, with the rise of *generative AI agents*, user queries tend to span multiple data sources. AI agents need consistent, high-level views of data to provide meaningful answers. A robust on-demand schema migration or virtualization strategy can streamline these queries without forcing constant data duplication.

### 1.1  Related Work

Category theory has influenced database theory for decades. Functorial data migration [1] represents schema changes as functors between categories, while topos-theoretic approaches [2] provide a unifying framework for logical interpretation. This paper extends these ideas with an emphasis on *real-time, on-demand transformations* and *AI-driven multi-model queries*.

### 1.2  Contributions

- We demonstrate the **usefulness of schema migrations** as an ongoing, rather than one-time, process.

- We illustrate **on-demand transformations** for queries crossing data models in real-world deployments.

- We detail how **generative AI agents** can exploit these transformations, discussing scenarios where data is scattered across heterogeneous backends.

- We provide **mathematical equations** and definitions capturing topos-theoretic data modeling.

- We conclude with a **sample Haskell code** outlining an implementation blueprint.

# 2 Usefulness of Schema Migrations in Multi-Model Contexts

**Schema migrations** enable:

1. **Evolving Business Requirements**: As organizations refine data requirements, schemas must adapt to new fields or constraints.

2. **Maintain Data Integrity**: Automated transformations preserve referential integrity and data consistency across changing structures.

3. **Heterogeneous Store Integration**: Relational, NoSQL, and graph stores have diverse schema representations. Migrations unify them under consistent, interpretable forms.

4. **Shorter Development Cycles**: Developers can roll out schema changes or expansions more quickly when migrations are automated and versioned.

# 3 Mathematical Foundations for Schema Transformations

Here, we briefly outline the category-theoretic constructs relevant to schema migration, illustrating how on-demand transformations can be formalized.

## 3.1 Schema as a Category

Let each schema be represented as a category $\mathcal{C}$, whose:

- **Objects**: Represent tables, collections, or entity types.

- **Morphisms**: Represent relationships (e.g., foreign keys, edges, or references).

A **schema instance** is then a functor $I : \mathcal{C} \to \mathbf{Set}$, mapping each object to a set of rows (or documents, or vertices) and each morphism to a function between sets.

## 3.2 Migrations as Functors

A schema transformation from $\mathcal{C}$ to another schema $\mathcal{D}$ is captured by a functor:

$$F : \mathcal{C} \longrightarrow \mathcal{D}. \tag{1}$$

Under such a functor, objects $c \in \mathcal{C}$ map to $F(c) \in \mathcal{D}$, and morphisms $c_1 \to c_2$ map to $F(c_1) \to F(c_2)$. At the level of *instances*, a *migrated* instance $I' : \mathcal{D} \to \mathbf{Set}$ can be composed with $F$ to produce:

$$I' \circ F : \mathcal{C} \to \mathbf{Set}. \tag{2}$$

This composition expresses how data in $\mathcal{C}$ is "reinterpreted" in $\mathcal{D}$'s schema.

### 3.3 Topos-Theoretic View

When dealing with complex schemas or partial knowledge, we may treat $\mathcal{C}$ and $\mathcal{D}$ as *sites* with Grothendieck topologies. The topos of sheaves (or presheaves) over these sites captures local data pieces that "glue" together globally [2].

In many real-world scenarios, we might define a *geometric morphism*:

$$\mathbf{Sh}(\mathcal{C}) \longrightarrow \mathbf{Sh}(\mathcal{D}),$$

where each sheaf in $\mathbf{Sh}(\mathcal{C})$ (data organized under schema $\mathcal{C}$) can be pulled back/pushed forward into $\mathbf{Sh}(\mathcal{D})$. This offers a powerful unifying framework for bridging data across multiple schemas and data models.

### 3.4 On-Demand Query Transformations

Instead of a single, one-time migration, let us consider a query $Q$ that references data from schema $\mathcal{C}$ but is formulated in terms of $\mathcal{D}$. Define a *schema translation functor* $F : \mathcal{C} \to \mathcal{D}$ plus a function $\phi$ that rewrites $Q_\mathcal{D}$ (a query in the language of $\mathcal{D}$) into a query $Q_\mathcal{C}$ if data physically resides under $\mathcal{C}$. Concretely, we can define:

$$\phi(Q_\mathcal{D}) = Q_\mathcal{C} \quad \text{such that} \quad \forall I : \mathcal{C} \to \mathbf{Set}, \quad Q_\mathcal{C}(I) \simeq Q_\mathcal{D}(I' \circ F). \tag{3}$$

This equation states that applying the rewritten query $Q_\mathcal{C}$ to instance $I$ yields the same "answer" as applying $Q_\mathcal{D}$ to the migrated instance $I' \circ F$, ensuring the user sees uniform results.

## 4 On-Demand Transformations & Generative AI Agents

### 4.1 Motivation

Generative AI agents often attempt to unify data from many sources to answer high-level questions. For example, an agent might receive a user query like: *"Find the total sales of Product A across all channels, including relational orders, NoSQL logs, and graph-based social recommendations."*

### 4.2 Multi-Model Setup

- **Relational Store** $(\mathcal{C}_R)$: Contains transactions and structured relationships.

- **Document Store** $(\mathcal{C}_D)$: Holds logs or aggregated JSON records.

- **Graph DB** $(\mathcal{C}_G)$: Stores user-product relationships or social network data.

An agent might parse the user request into a "universal" high-level query in $\mathcal{D}$. If the needed data is physically split across $\mathcal{C}_R, \mathcal{C}_D, \mathcal{C}_G$, the system either:

1. Defines transformations $F_R, F_D, F_G$ to unify them in a single schema $\mathcal{D}$, or

2. Dynamically rewrites subqueries for each store and aggregates the results.

### 4.3 Importance of Automated Migrations

- **Minimal Data Duplication**: On-demand transformations can pull or transform data at query time, negating the need for perpetual synchronization.

- **Consistent Data Models**: By harnessing formal functors and geometric morphisms, organizations ensure that the partial or local data automatically aligns with the global schema constraints.

- **Generative AI Synergy**: AI agents can learn or infer schema mappings from partial examples, and rely on validated category-theoretic transformations for reliable, real-time integration.

## 5 Benefits and Practical Observations

### 5.1 Reduces Operational Risk

Traditional "big-bang" schema migrations can be risky and time-consuming. On-demand transformations localize changes to the relevant queries or data subsets, decreasing downtime.

### 5.2 Enables Incremental Adoption

Organizations can adopt a unified approach *gradually*, layering new transformations only when necessary. This fosters agility and lowers the barrier to rolling out new data-driven features.

### 5.3 Accelerates AI-Driven Insights

Generative AI agents benefit from consistent, on-demand transformations because:

1. They can operate on a "universal schema" while the underlying data remains in place.

2. They can quickly adapt to changes in the physical data layout.

3. They avoid manually-coded ETL pipelines and harness logic-based transformations for consistent results.

## 6 Sample Haskell Implementation

In practice, a functional language like Haskell can elegantly express these transformations. We present a separate-file code snippet in Section A that demonstrates how categories, functors, and on-demand transformations might be modeled.

## 7 Conclusion

We have shown how schema migrations, traditionally seen as a one-time event, can instead be deployed *on demand*, enabling queries and generative AI agents to seamlessly traverse multi-model data. By leveraging the formalism of category and topos theory, we achieve consistency, modularity, and extensibility. Future work involves deeper integration of proof assistants for verifying correctness, real-time streaming transformations, and further AI-driven schema discovery.

# References

[1] D. I. Spivak. *Functorial Data Migration.* arXiv:1009.1166, 2010.

[2] P. T. Johnstone. *Sketches of an Elephant: A Topos Theory Compendium.* Oxford University Press, 2002.

# A   Sample Haskell Code

See file: `onDemandMigration.hs` for a minimal example of representing schemas, morphisms, and data transformations using Haskell.