# SAMARTH

- Designed, Developed and Implemented by IIC, University of Delhi
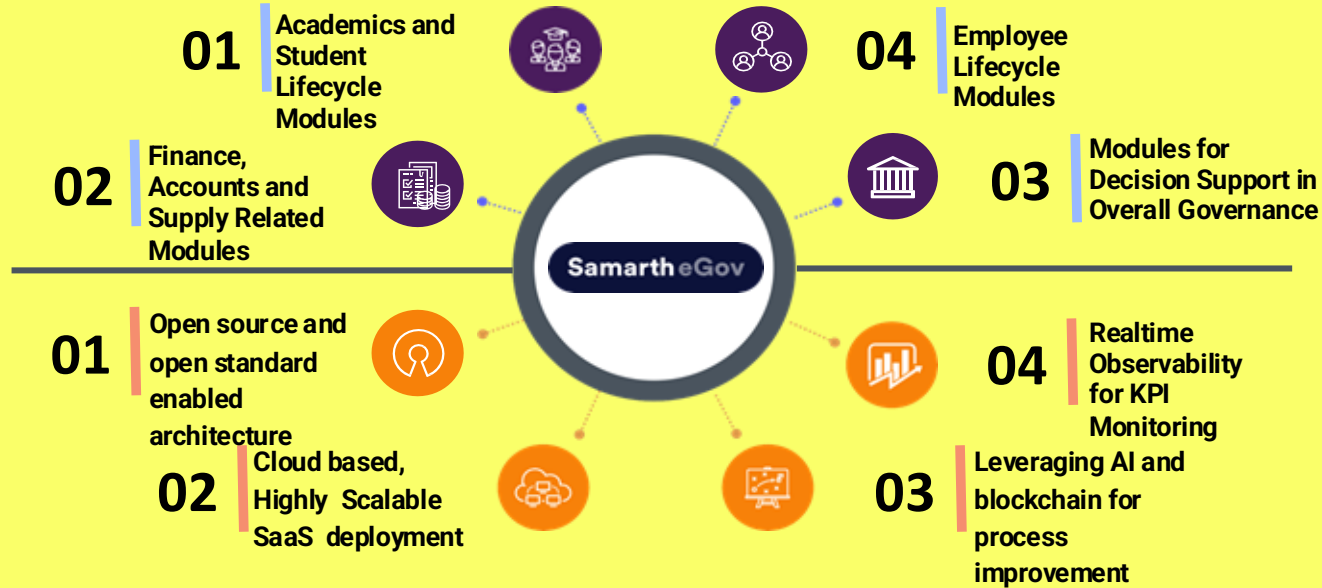
## SAMARTH's USP

Transform your institution with SAMARTH eGov - a powerful, scalable, and future-ready

platform built for higher education excellence.

**Samarth** eGov

**Building Future Ready Digital Campuses**

ClickHouse Gurgaon Meetup

Samarth eGov

# Introduction

Purpose built platform for Higher Education Institutions (HEIs) to deploy e-governance system for administration and management of higher education processes.

**01** | Academics and Student Lifecycle Modules

**04** | Employee Lifecycle Modules

**02** | Finance, Accounts and Supply Related Modules

**03** | Modules for Decision Support in Overall Governance

SamartheGov

**01** | Open source and open standard enabled architecture

**04** | Realtime Observability for KPI Monitoring

**02** | Cloud based, Highly Scalable SaaS deployment

**03** | Leveraging AI and blockchain for process improvement

SAMARTH emerges as the point of convergence for higher education policies, compliances and services.

SamartheGov

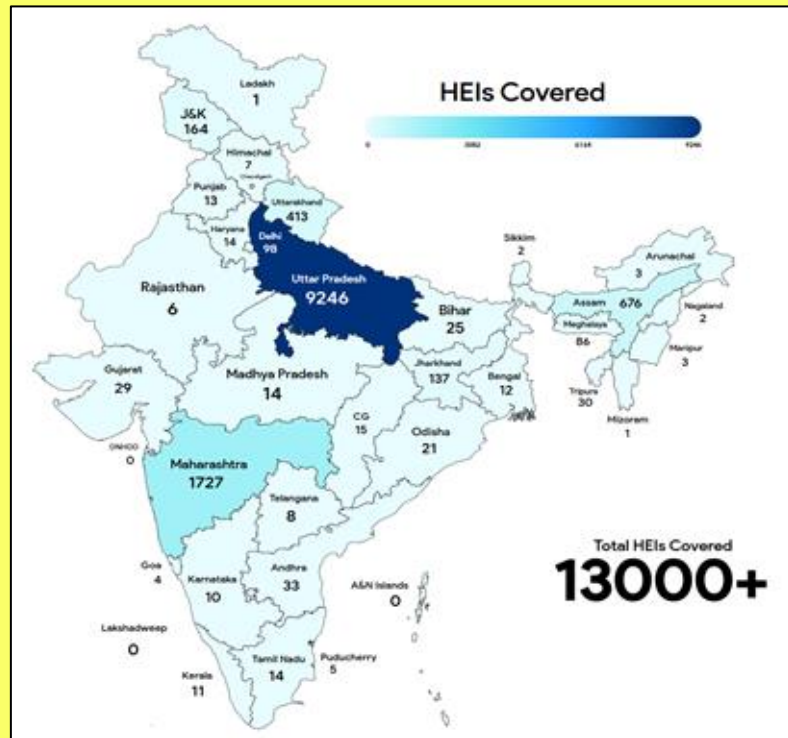# SAMARTH IMPLEMENTATION STATUS - SCALE & SPREAD

Connected **13000+ HEIs** across **32 States/UTs** and **435 Districts/Cities**

**STUDENT LIFECYCLE – FROM ADMISSIONS TO CONVOCATION**

**EMPLOYEE LIFECYCLE – FROM RECRUITMENT TO RETIREMENT**



HEIs Covered

Ladakh 1
J&K 164
Himachal 7
Punjab 13
Haryana 14
Uttarakhand 413
Delhi 98
Rajasthan 6
Uttar Pradesh 9246
Sikkim 2
Arunachal
Assam 676
Nagaland
Meghalaya 86
Manipur 3
Bihar 25
Jharkhand 137
Bengal 12
Tripura 30
Mizoram 1
Gujarat 29
Madhya Pradesh 14
CG 15
Odisha 21
DNH DD
Maharashtra 1727
Telangana 8
Goa 4
Karnataka 10
Andhra 33
A&N Islands
Lakshadweep 0
Kerala 11
Tamil Nadu 14
Puducherry 5

Total HEIs Covered
**13000+**

**Samarth eGov**

**ClickHouse Gurgaon Meetup**

# Understanding the SAMARTH Data Landscape

**21.7 M**
Records in Student Registry

**254 K**
Records in Employee Registry

**1.24 M**
Payroll Statements Processed Online

**966 K**
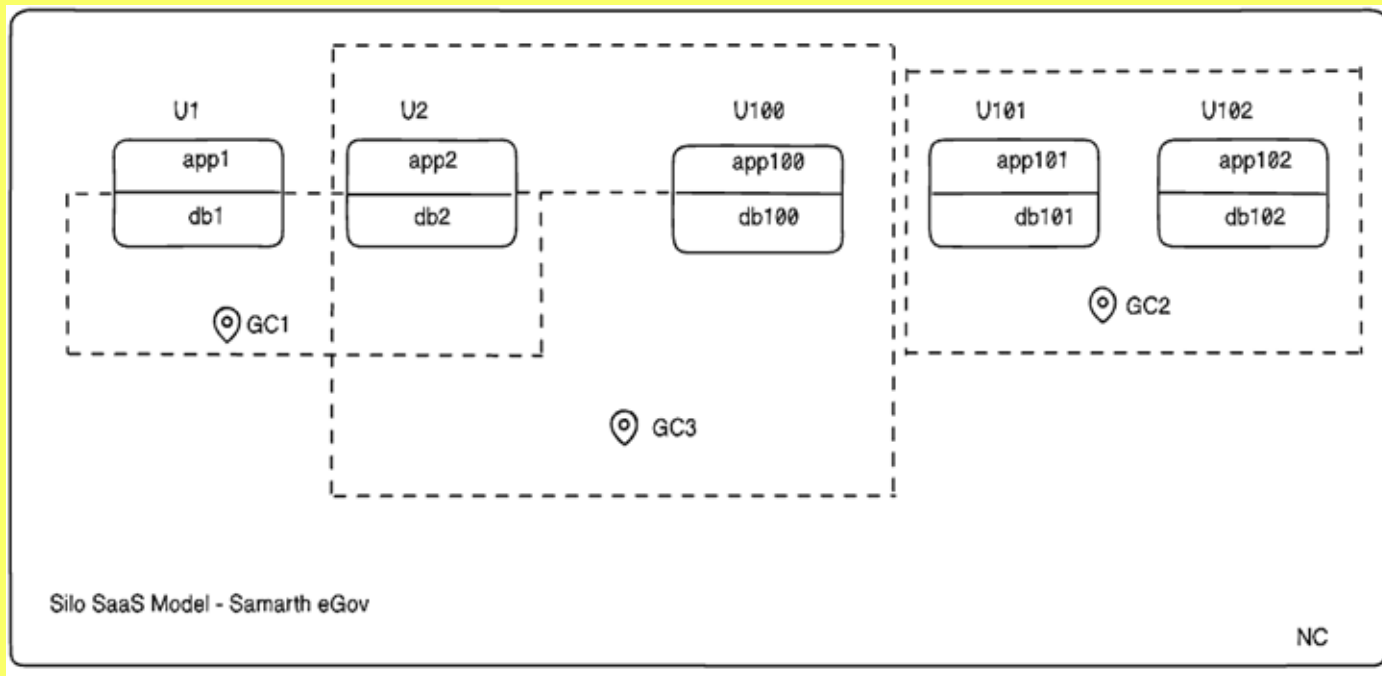Active Students

**228 K**
Active Employees

**1.01 M**
Leaves Processed Online

- Analytical datasets operating at million-level record scale
- Continuous, high-volume write amplification workloads

Samarth eGov

**ClickHouse Gurgaon Meetup**

# Dimensions That Matter!

- Nation Top Level

- Geographic hierarchies (state, district, region)

- Institutional and organizational hierarchies (university, college, department)

- Temporal dimensions (academic year, financial year, seasons, time)

- Staff Attributes (designation, employment nature, service status)

- Programme and disciplin... ... ... (stream, course, specializations, elective)

- Demographic attributes (category, reservation, gender)

.. etc

ClickHouse Gurgaon Meetup

# The Dimensional Multiplier: Why Complexity Scales ?



Silo SaaS Model - Samarth eGov

NC - National Cluster
GC* - Geographic Cluster
U* - Individual Subscriber

# Symptoms

DB Schemas can vary - Heterogeneous

Computational Latency due to High Dimensionality

Performance Bottlenecks in Deep-Dive Analysis

Excessive CPU and memory consumption due to complex transformations and joins

Heavy on-the-fly aggregations leading to slower response times and unpredictable performance

Samarth eGov

# Journey @ Samarth Analytics Platform

- **Gen1 – Metric-Driven ETL with self defined KPI templates for each subscriber. Additional custom dashboards for NC/GC aggregations**
  - *Technology –* `MySQL + Grafana`
  - *Issues Faced –*
    - Rigid schema, no ad-hoc analysis or correlation
    - Redesign required for every new metric
    - Access provisioning to every end users
    - Cross                 Subscriber                 aggregations

**ClickHouse Gurgaon Meetup**

Samarth eGov

# Journey @ Samarth Analytics Platform

- **Attempt 1 for Gen2(.v1) — Moving to Data Warehouse**

  - *Technology —* `Redshift + Quicksight`
  - *Issues Faced —*
    - Schema mismatches and modeling errors
    - Data Ingestion overhead
    - High operational effort and cost

*Could not materialize to production*

# Journey @ Samarth Analytics Platform

- **Attempt 2 for Gen2(.v2) – Data Warehouse using CDC (ELT)**

  - *Technology –* `CDC + Clickhouse + Superset`   
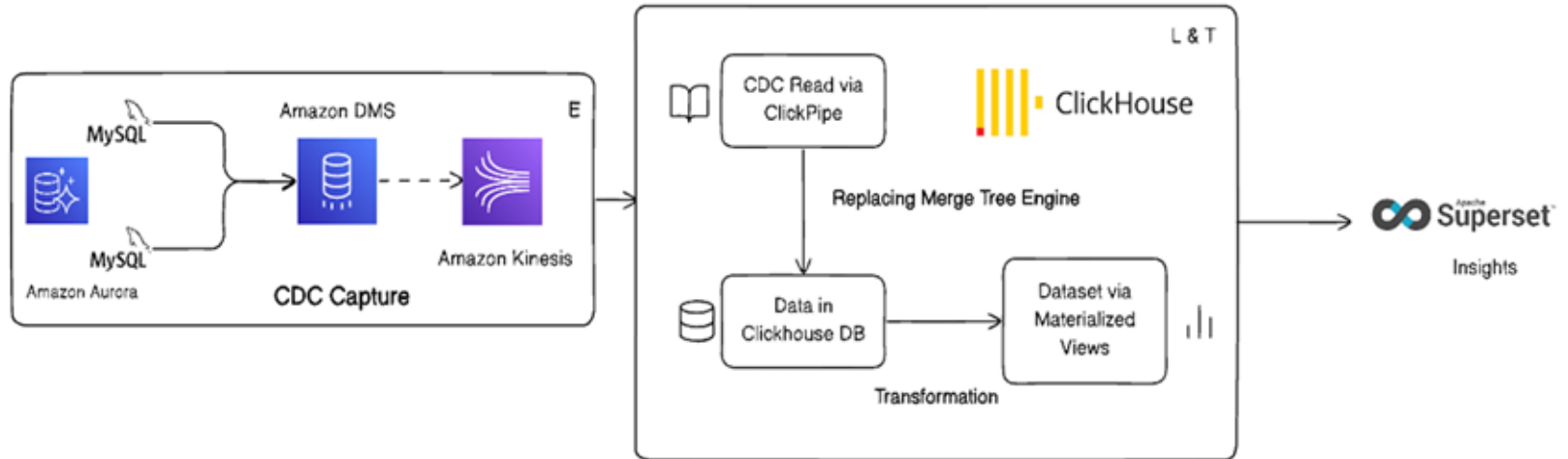
  - *Facts –*

    - Schema-Flexible by Design

    - Reduced Pipeline Complexity via Click Pipes

    - Simplifies architecture and lowers operational cost

    - High performance for group aggregations and analytics

    - Fast Execution by using refreshable materialized views

Samarth eGov

# Architecture

End-to-end CDC ingestion via *Clickpipes* into ClickHouse, utilizing the **ReplacingMergeTree** engine and **Refreshable Materialized Views** to deliver *near real-time analytics*



ClickHouse Gurgaon Meetup

# ClickHouse transformed high-dimensional, high-velocity data into low-latency, cost-efficient analytics at scale

Samarth eGov

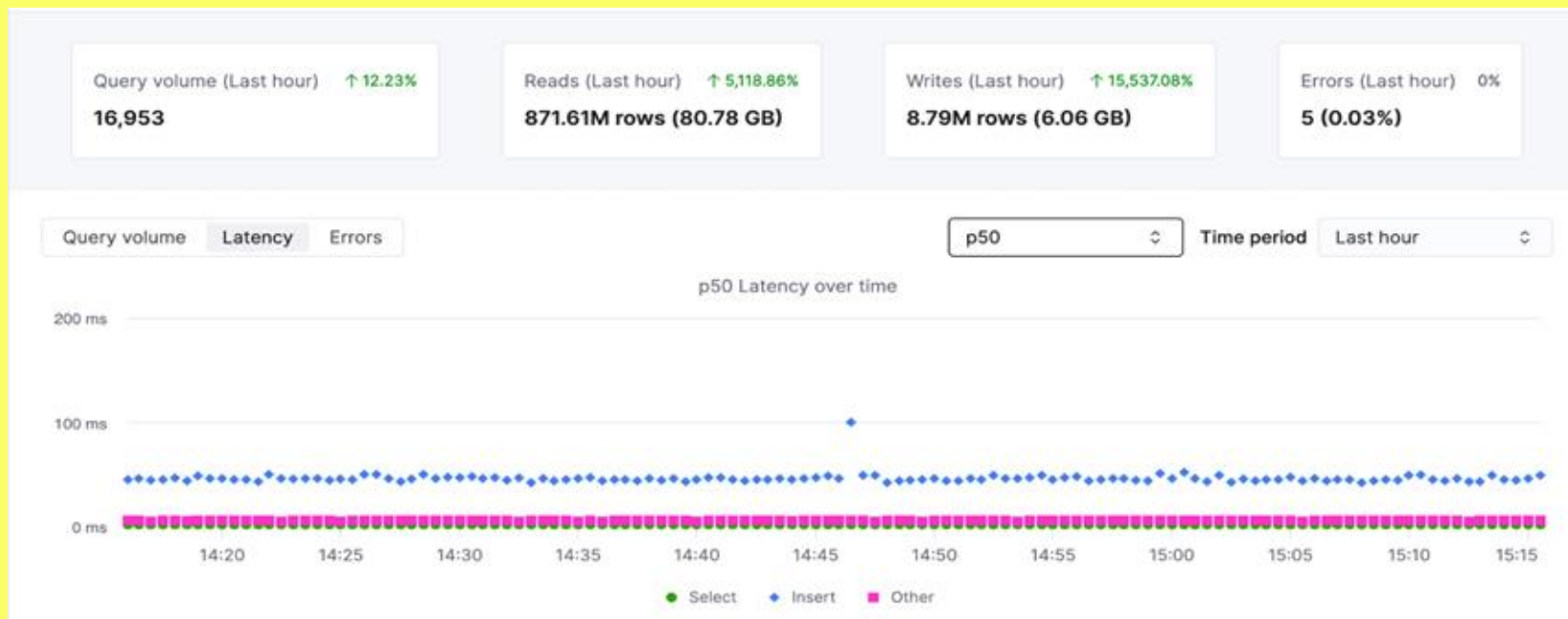| Aspect | v1 – MySQL + Grafana (Classic ETL) | v1.5 – Redshift + QuickSight | v2 – ClickHouse + CDC + Superset |
|---|---|---|---|
| Architecture | Direct queries on source DB with ETL | Centralized warehouse + BI | Streaming CDC-based analytics |
| Ingestion | Batch ETL | Heavy batch loads | Continuous CDC loading |
| Schema Handling | Rigid, predefined metrics | Strict schema, load failures on mismatch | Schema-flexible, non-blocking ingestion |
| Data Model | Siloed, metric-driven | Silo-oriented tables | Unified, analytics-oriented |
| Failure Impact | Query failures affect dashboards | Load failures block process | Errors isolated in separate tables |
| Correlation Across Data | Not supported | | Fully supported |
| Query Performance | Slow at scale | | Fast group aggregations and analytics |
| Real-time Capability | No | | Near real-time |
| Visualization Layer | Grafana | QuickSight | Superset |
| Scalability | Poor | | High (designed for scale) |

Samarth eGov

After heavy aggregations and complex fact-dimension joins, e.g. -

Duration Taken: 24.7 seconds

Read Data: 4.85 GB | Written Data: 6.02 GB

Total Rows Materialized: 28.7 million

Samarth eGov

# After heavy aggregations and complex fact-dimension joins, e.g. -

# Thank You!

→ **Kunal Sharma**

**Associate, Data Engineering and AI**