# Einblick

**Product Overview for Databricks**

# AI will create 7 billion new data scientists

Historically, data analysis was a complex technical endeavor. Today, AI empowers everyone to find the right answers simply by asking natural language questions.



**1940**



**1960**



**2020**



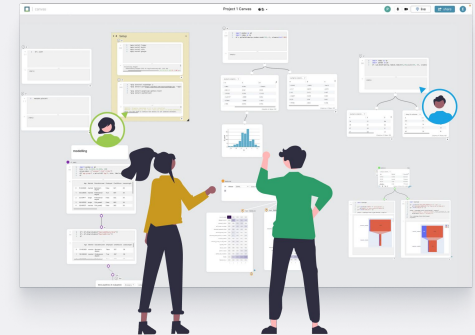Tell me how product ratings affect sales?

**Now**
Natural language replaces programming

Einblick

# Current Trends

## 98% of Python is written in local editors



| | |
|---|---|
| 36% | VS Code |
| 29% | PyCharm |
| 21% | Other non-sharable |
| 7% | Vim |
| 5% | Jupyter Notebook |
| 2% | Jupyter Lab |

## BI / analytics and data science are converging



## Democratization of data analytics / science

# What is Einblick?

A next-generation, AI-native, multi-modal data notebook to build workflows and data apps



## Natural language

Explore data, build workflows and refine cells by having a conversation with the system.

## No-code

No-code operators for common tasks (visualizations, table manipulations, automated ML, etc.) to enhance productivity and empower users of all technical backgrounds. Input controls to build data apps.

## Code

Native support for SQL and Python. Pull data from BiqQuery, Snowflake, etc. and pipe into Python cells which are functionally equivalent to Jupyter notebooks.

Einblick

4

# Target Persona

## Editors

Data practitioners, folks that are familiar with Python and / or SQL
*(still early, but clear indication that barrier to entry significantly lowered through AI)*

## Viewers

Citizen data scientists, stakeholders in the analytics process, oftentimes less technical

Einblick

# Use Cases

## Advanced reporting



Data in SQL & wrangling in Python
Recurring

## Data science experimentation



Python, ML, SQL, no-code
Ad-hoc

# Who are we?

Making data analytics a more effortless, efficient, and collaborative experience.
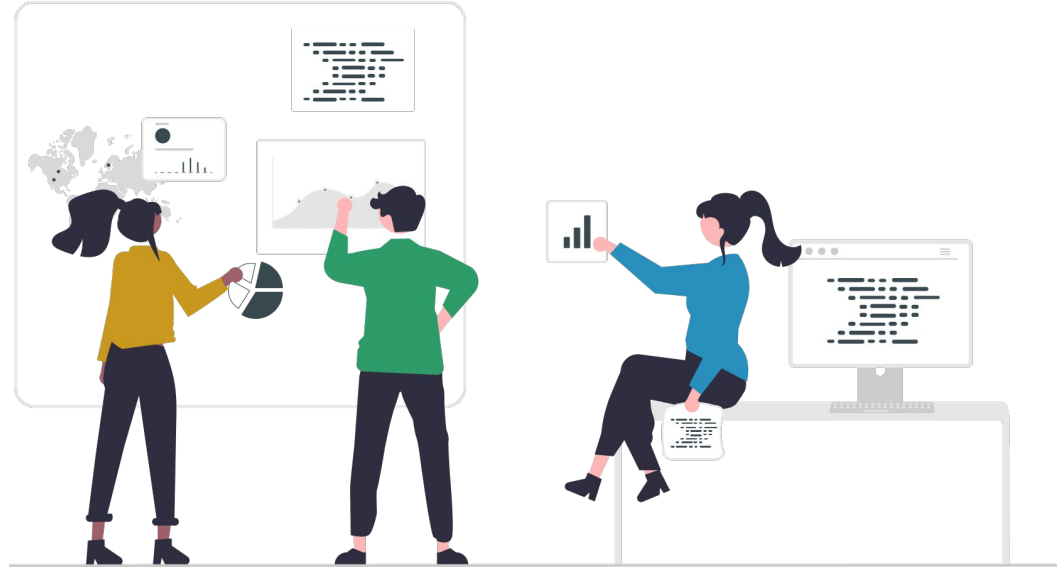
Founded in 2020 as a remote first company based primarily in Boston and NYC with talent in applied machine learning and frontend engineering

6+ years of research at the intersection of human-computer interaction (HCI) and machine learning at MIT and Brown University.

Amplify PARTNERS | DARPA | DELL | FLYBRIDGE CAPITAL PARTNERS | intel capital | SAMSUNG NEXT

10M in total funding

Einblick

## Emanuel Zgraggen
CEO / Co-founder

- Ex-Postdoc @ MIT
- Ph.D. in Computer Science from Brown University
- Passion for products at the intersection of HCI and machine learning / AI
- Research[1] in HCI and applied ML, building interactive tools for data analysis and data science

## Philipp Eichmann
CXO / Co-founder

- Ph.D. in Computer Science from Brown University
- Life-long UI/UX enthusiast and engineer
- Passionate about designing and building tools for data science
- Conducted research[2] in HCI to democratize access to data science

[1] bit.ly/3vKvNO2
[2] bit.ly/3Oh0gcO

Einblick

8

# User Interface

Einblick

# User Interface



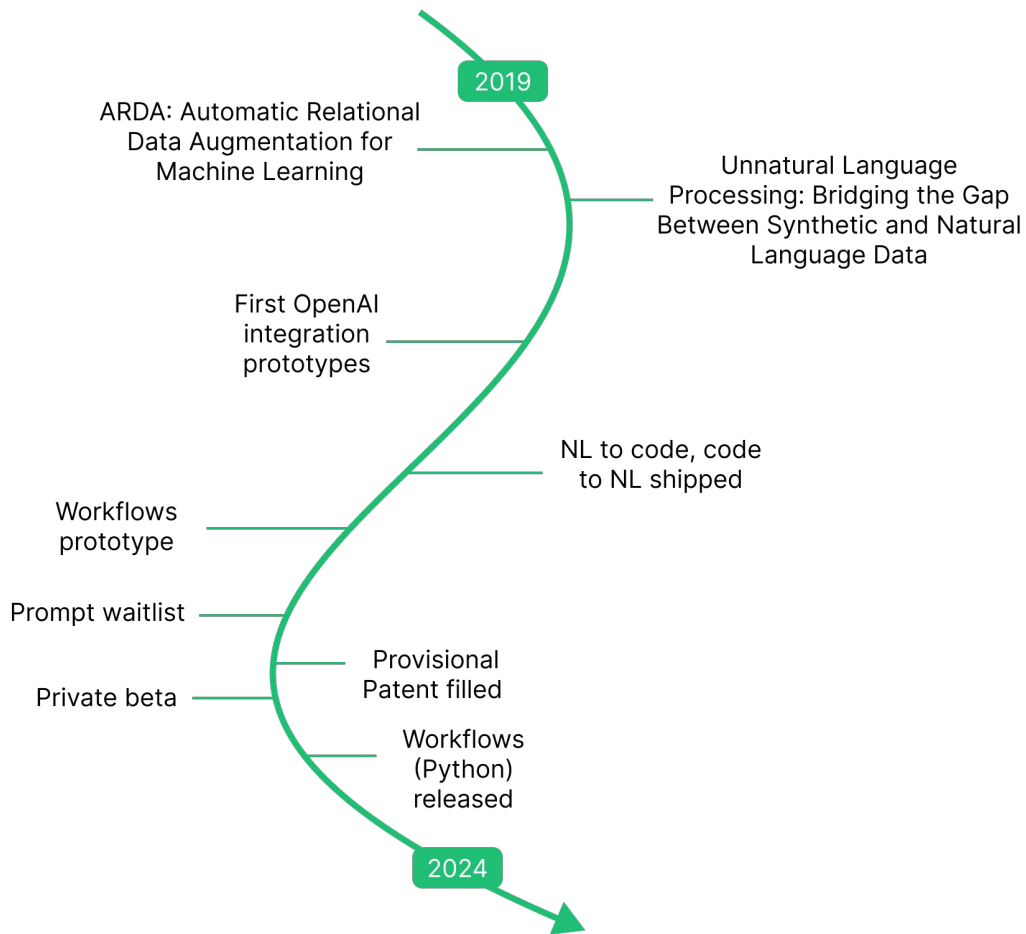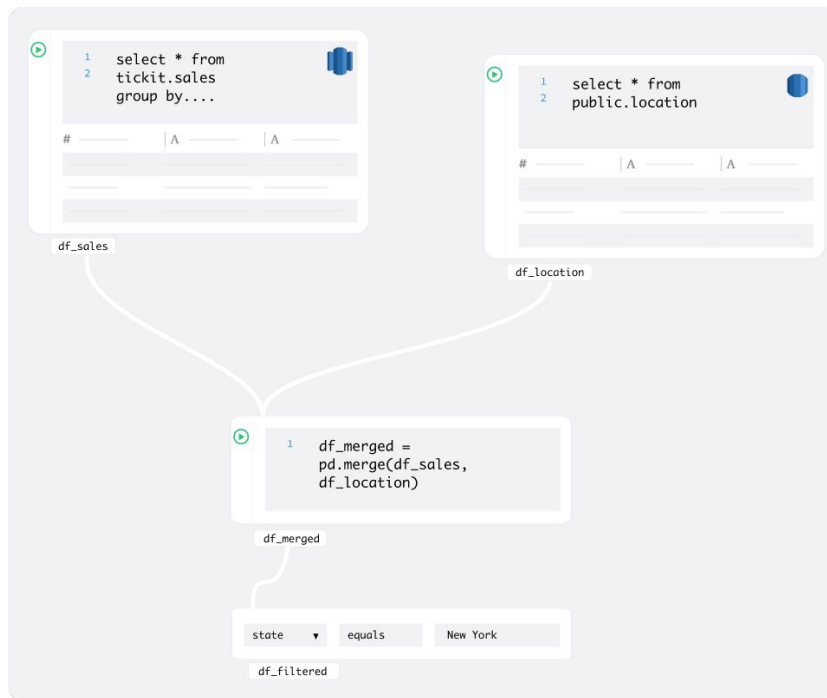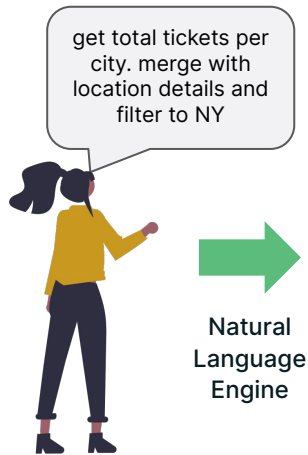| Problem | Solution |
|---|---|
| Notebooks preclude non-technical collaborators | Completely re-thought notebook/canvas design, library of no-code cells |
| Collaboration is crucial | Collaboration is a first principle: canvas as an emerging means for collaboration, real-time video/audio streaming |
| The possibilities for no-code operators are endless | Extensible architecture and plugin Infrastructure |
| BI and Data Science are converging | Multi-modal environment: no-code, code, and natural language |
| Notebooks cannot easily be turned interactive, shareable experiences | User controls and reactive execution, data apps |
| Interpretability and reproducibility of notebooks can be confusing and hard | Automatic dependency resolution using static code analysis, exposed to the user |

Einblick

# Natural Language Engine

# From text to data workflow

**Why we care:**

- **Time saver for experts**
- **Lower barrier of entry for non-experts**

**2019**

ARDA: Automatic Relational Data Augmentation for Machine Learning

Unnatural Language Processing: Bridging the Gap Between Synthetic and Natural Language Data

First OpenAI integration prototypes

NL to code, code to NL shipped

Workflows prototype

Prompt waitlist

Provisional Patent filled

Private beta

Workflows (Python) released

**2024**

**Targeted**: best possible AI engine for the data analytics and data science

**Multimodal**: mix no-code, Python, SQL

**Extensible**: incorporate new task types and modalities

**Adaptable**: keep up to date with the most recent advancements

# Example: "Plot age as a histogram"



**GitHub Copilot** ❌

Picked the second out of 5 suggestions. Executing code leads to errors:
- Wrong dataframe name (df instead of marketing)
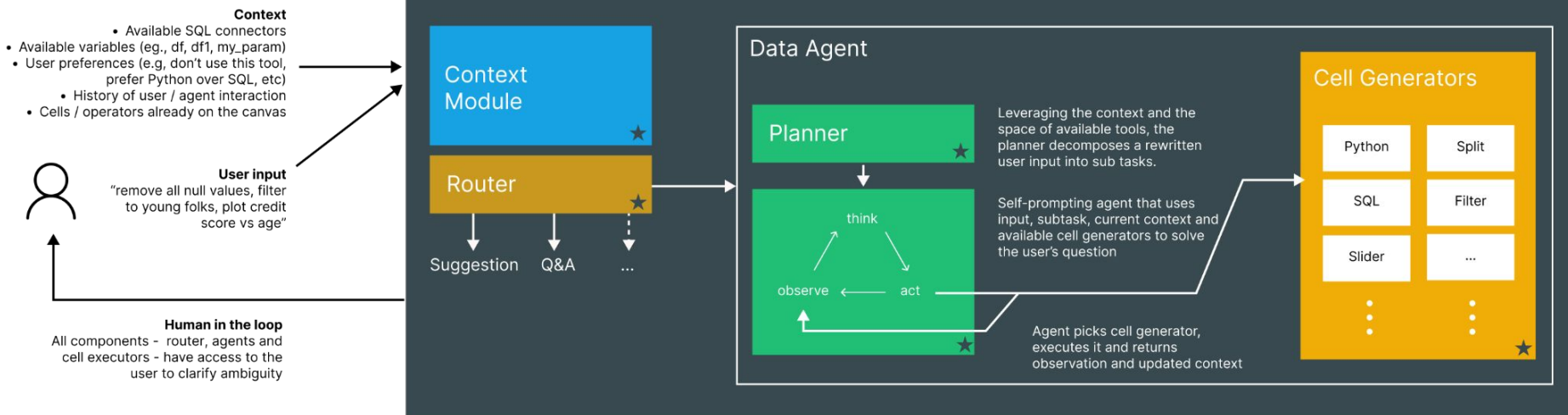- Wrong column name (age instead of Age)

**Einblick** ✅

# The hard stuff nobody talks about

| Problem | Solution |
|---|---|
| LLMs are not a product | Full infrastructure is required, connectors, no-code plugins, Python / SQL execution |
| Latency latency latency! | Various UI techniques<br><br>Adaptive routing based on context |
| Hallucinations | Give LLMs a way out<br><br>Human in the loop |
| Testing / Logging<br>- Small changes have big impact<br>- Standard engineering techniques fail | Set up extensive infrastructure<br><br>Automated regression testing using LLMs for comparison |
| Robustness<br>- LLMs will return garbage at times | Schema validation, code parsing<br><br>Self-correction |
| Multi-modality<br>- LLMs get confused when mixing languages | Built-in abstraction layer to translate everything in context (e.g., SQL, no-code) to Python |
| Context<br>- Limited window size<br>- Tokens are expensive and slow<br>- Recency bias | Context module with static and dynamic part<br><br>Multi hierarchy vector stores<br><br>Many UI based heuristics |

# Natural Language Engine Architecture

**Context**
- Available SQL connectors
- Available variables (eg., df, df1, my_param)
- User preferences (e.g, don't use this tool, prefer Python over SQL, etc)
- History of user / agent interaction
- Cells / operators already on the canvas

**User input**
"remove all null values, filter to young folks, plot credit score vs age"

**Human in the loop**
All components - router, agents and cell executors - have access to the user to clarify ambiguity

### Context Module

### Router

Suggestion    Q&A    ...

## Data Agent

### Planner

Leveraging the context and the space of available tools, the planner decomposes a rewritten user input into sub tasks.

think

observe ← act

Self-prompting agent that uses input, subtask, current context and available cell generators to solve the user's question

Agent picks cell generator, executes it and returns observation and updated context

### Cell Generators

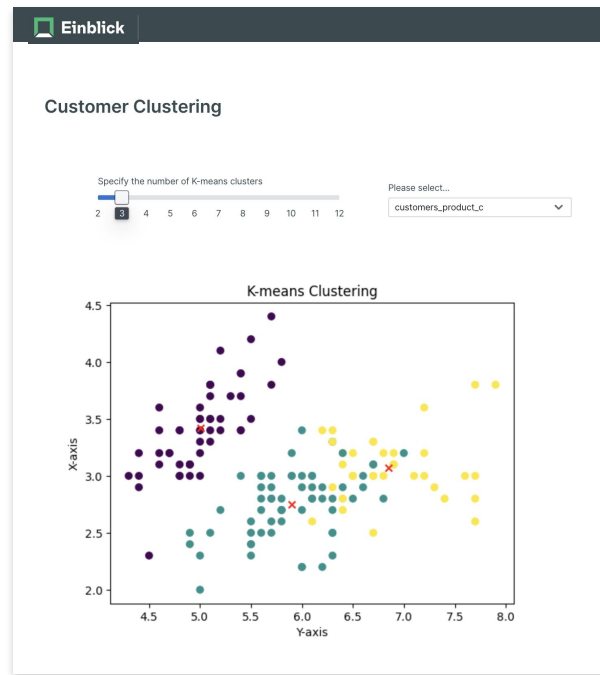| Python | Split |
|--------|-------|
| SQL | Filter |
| Slider | ... |

★ **LLM Powered**
- All major components within the system are LLM powered
- We currently use OpenAI's gpt-3.5-turbo, however architecture is model agnostic

Einblick

# Data Apps & Plugins

Einblick

# Data Apps



Anything on a canvas can be turned into an **interactive, shareable data app,** with a few clicks. **No code required.**

# Plugins



**Build or generate** plugins **without code**, on the canvas. Plugins are **available as no-code cells** through the menu, and can be made **accessible to the Natural Language Engine**.

# www.einblick.ai