

# Epigene Labs at JOBIM 2025: Data Engineering, Open Science and Scalable Tools in Action

August 4, 2025

From July 8th to 11th, three Epigeners traveled to Bordeaux to take part in [JOBIM 2025](#), a key conference for computational biology. Epigene Labs presented two posters and gave a live demo of its open-source tool, InMoose. Here's a look at what we shared with the scientific community.

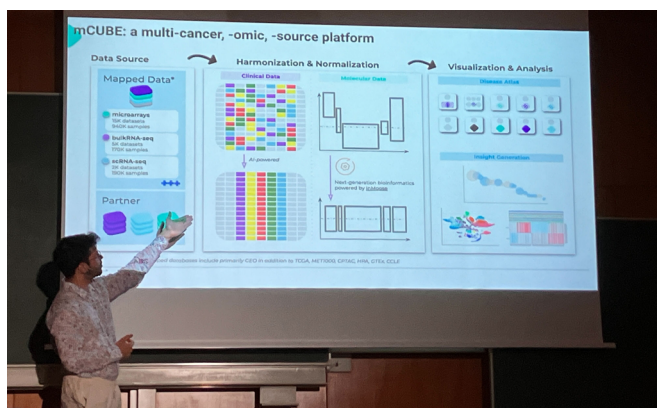
## InMoose: from a bold bet to a robust bioinformatics library

A few years ago, Epigene Labs made a bold move: simplifying its **data science** stack by migrating fully to Python – leaving behind the traditional R/Python dual setup used in most bioinformatics workflows. This decision meant porting essential R-based tools to Python, as no equivalents were available at the time – an ambitious **data engineering challenge**.

The effort paid off. This initiative led to InMoose, a Python-native library that now reimplements trusted tools for bulk transcriptomic analysis, including limma, edgeR, DESeq2, ComBat, and ComBat-Seq. With three peer-reviewed validations, InMoose proves that robust **data engineering** can improve reproducibility without compromising performance.

On July 8th, Maximilien Colange, lead developer of InMoose, walked attendees through the origins of the project and demonstrated how to use it effectively. Several presentations during the conference also featured research conducted using the InMoose environment, a strong sign of its growing adoption by the community.

Curious to try it yourself? You can access the tutorial as a [shared notebook](#).



# Epigene Labs at JOBIM 2025: Data Engineering, Open Science and Scalable Tools in Action

August 4, 2025

## Cohort sizing: smarter data engineering for scalable science

Robust scientific findings depend on reproducibility and generalizability. Larger cohorts are key to capturing biological variability and producing trustworthy results. Yet, omics data often suffers from fragmentation – different technologies, protocols, and naming conventions make datasets hard to combine, despite the clear value of doing so.

On July 10th, Maximilien Colange presented a methodology that helps quantify the trade-off between cohort size and analytical accuracy. This work provides a concrete answer to a long-standing question: How much value are we losing by working with small, fragmented datasets?

This work advocates for making **data integration** a strategic, well-funded **data engineering** priority, not an afterthought.

To support the community, we've also released an [online calculator](#) to help researchers assess the impact of their cohort size.

### A priori Estimation of Reproducibility Odds Informs the Sizing of Omic Data Cohorts

M. Colange<sup>1</sup>, A. Nordori<sup>1</sup>, and A. Behdenna<sup>1</sup> JOBIM 2025  
Bordeaux  
<sup>1</sup>Epigene Labs, Paris, France

#### Introduction

- Fragmented Data Hinders Progress:** Omic studies are constrained by fragment or poorly integrated datasets, weakening generalizability and reproducibility.
- Data Integration is Under-Resourced:** Data integration is typically approached on a "best effort" basis, with little guidance on how much effort or resources are truly needed.
- Costs Remain Invisible:** The scientific and opportunity costs of limited integration are widely acknowledged but rarely quantified in a rigorous and systematic way.

#### Contribution

- Quantifying What's at Stake:** We introduce mathematical formulas that link cohort size and statistical power, making the cost of limited integration explicit.
- Practical Tools for Study Design:** These ready-to-use formulas apply broadly across data types and can inform both new and secondary analyses.
- Enabling Strategic Commitment:** By revealing the price of underpowered studies, we aim to shift data integration from an ad-hoc task to a justified, well-resourced priority.

#### Quantifying the gap between observations and the hidden signal

**Law of large numbers**  
The distribution of observations converges to the true distribution. "the more samples, the better"

**Dvoretzky-Kiefer-Wolfowitz bound**  
Gives a rate of convergence for the law of large numbers. "n samples needed for a 95% CI on the signal distribution narrower than  $\epsilon$ "

n	$\epsilon$
500	6.07%
1,000	4.29%
5,000	1.92%
10,000	1.36%

#### Linking CI width, confidence level and cohort size

$$1 - \alpha \leq 2ke^{-2n\epsilon^2}$$

confidence level on the approximation of the signal distribution:  $1 - \alpha$   
number of observations:  $n$   
number of dimensions:  $k$   
CI width:  $\epsilon$

Question	Answer
How many samples to have 95% CI narrower than 5%? ( $\alpha = 95\%$ , $\epsilon = 5\%$ )	<b>N = 2,441</b>
Width of the 95% CI with 2000 observations? ( $\alpha = 95\%$ , $n = 2,000$ )	<b><math>\epsilon = 5,5\%</math></b>
What CI are narrower than 5% with 2,000 observations? ( $\epsilon = 5\%$ , $n = 2,000$ )	<b><math>\alpha = 54,6\%</math></b>

#### Application to cancer datasets from GEO

**Microarray** (for de novo data generation):  $k=5,000$  features. Target cohort size ~ 2,500 samples for 95% CI narrower than 5%.  
N = 2,500,  $\epsilon = 5\%$  (16.7% of available data)  
N = 15,000,  $\epsilon = 2\%$  (6x more precise)  
N = 38,  $\epsilon = 40\%$  (0.25% of available data)  
6x more precise, 65x more samples

**RNA-Seq** (for de novo data generation):  $k=5,000$  features. Target cohort size ~ 2,500 samples for 95% CI narrower than 5%.  
N = 2,500,  $\epsilon = 5\%$  (16.7% of available data)  
N = 15,000,  $\epsilon = 2\%$  (6x more precise)  
N = 32,  $\epsilon = 44\%$  (0.2% of available data)  
11x more precise, 78x more samples

#### Would you rather invest in under-powered de novo data, or in data integration capabilities?

- Data- and model-agnostic:** our formulas only depend on the number of features modeled. This guarantees their wide applicability.
- Fast estimates:** our formulas can be used to inform project feasibility, before experimental design is finalized and before data is collected.
- Untapped potential:** without data integration capabilities, public data remains under-utilized. Our study shows that only 15% of GEO data would improve study precision by a factor 8 to 11.
- Addressing biases in available data:** whatever the cohort size, representativity is key. Data integration allows tailor-made cohorts, built to alleviate biases.
- Challenges in data heterogeneity:** evolving disease classifications, non-standardized nomenclatures, improving sequencing technologies, changing gene name references...
- Solutions exist:** batch effect correction, gene name harmonization, AI-powered clinical metadata cleaning...

CHECK OUR ONLINE CALCULATOR [Contact](#) [Maximilien Colange, PhD](mailto:Maximilien.Colange@epigenelabs.com) [ajournal@epigenelabs.com](mailto:ajournal@epigenelabs.com) The authors have no conflict of interest to declare. CHECK OUT

# Epigene Labs at JOBIM 2025: Data Engineering, Open Science and Scalable Tools in Action

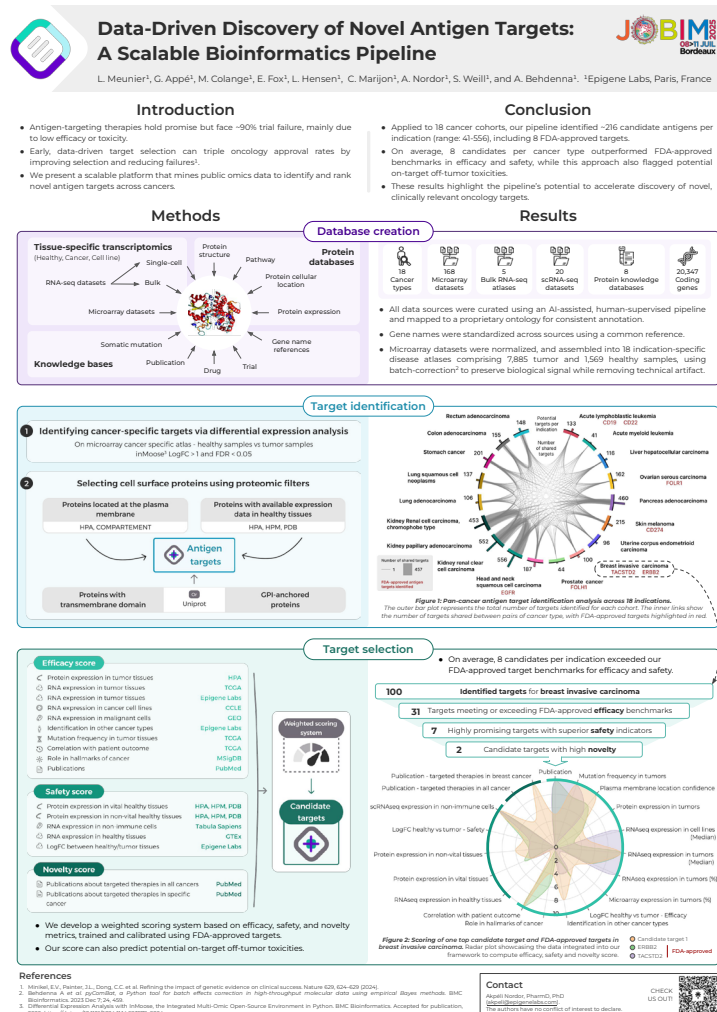
August 4, 2025

## Target discovery pipeline: smarter paths to cancer therapies

Antigen-targeting therapies have shown great promise in oncology – but most candidates still fail in clinical trials due to safety or effectiveness concerns. Identifying better targets earlier in the pipeline can dramatically improve success rates.

On July 10th, Léa Meunier, our senior computational biologist, presented Epigene's **scalable, data engineering-driven platform** that leverages public omics datasets to discover and rank new antigen targets across various cancer types.

The pipeline, validated across 18 cancer types, uses a human-in-the-loop AI approach for **data curation and harmonization** (e.g., normalization, batch effect correction, gene name standardization), to integrate a comprehensive compendium of data: 8 protein knowledge databases, 5 bulk RNA-seq atlases, 20 tissue-specific single-cell RNA-seq datasets and 168 microarray studies assembled into cancer-specific cohorts with over 9,400 samples and 20,000 genes.



# Epigene Labs at JOBIM 2025: Data Engineering, Open Science and Scalable Tools in Action

August 4, 2025

Key results:

- **216 candidate antigens** per cancer type on average
- **8 FDA-approved antigens** rediscovered through the pipeline
- **8 novel candidates per indication** that exceed FDA benchmarks for efficacy and safety

This work highlights how **data engineering in biopharma** can accelerate and de-risk oncology drug development.

*Authors: Maximilien Colange & Léa Meunier*