

# ClusterProfiler-LLM: NSCLC Case Study

Automated Interpretation of Nature Genetics (2025) Data

2026-02-06

## Table of contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Materials and Methods</b>	<b>2</b>
2.1	Data Source . . . . .	2
2.2	Setup . . . . .	2
<b>3</b>	<b>Analysis Workflow</b>	<b>2</b>
3.1	1. Data Preparation . . . . .	2
3.2	2. Context Construction . . . . .	3
3.3	3. Standard Enrichment Analysis . . . . .	4
<b>4</b>	<b>Results: LLM-Driven Interpretation</b>	<b>5</b>
4.1	Task 1: Cell Type Annotation (Hypothesis Verification) . . . . .	5
4.2	Task 2: Phenotypic Characterization . . . . .	6
4.3	Task 3: Mechanism Interpretation (Multi-Agent Deep Mode) . . . . .	7
4.4	Task 4: Hierarchical Interpretation . . . . .	8
<b>5</b>	<b>Conclusion</b>	<b>10</b>

## 1 Introduction

This document demonstrates the capability of `clusterProfiler-LLM` to reproduce and interpret key findings from a high-complexity single-cell study. We analyze the dataset from *Nature Genetics* (2025): “Multi-omic profiling highlights factors associated with resistance to immuno-chemotherapy in non-small-cell lung cancer”.

We compare the automated output of our multi-agent framework against the manual conclusions drawn by the original authors, specifically focusing on:

1. **Mechanism Interpretation:** Tumor cell states (e.g., NRF2-mediated stress response).
2. **Cell Type Annotation:** Identification of specific subsets (e.g., SPP1+ Macrophages).
3. **Phenotypic Characterization:** Immune states (e.g., TIGIT/CTLA4-driven exhaustion).

## **i** Note

Online      HTML      version:      [https://yulab-smu.top/clusterProfiler-LLM-Supplementary/NSCLC\\_Case\\_Study.html](https://yulab-smu.top/clusterProfiler-LLM-Supplementary/NSCLC_Case_Study.html)

## 2 Materials and Methods

### 2.1 Data Source

- **Study:** Yan et al., *Nature Genetics*, 2025.
- **Context:** NSCLC treated with neoadjuvant ICB + chemotherapy.
- **Key Original Findings:**
  - **SPP1+ Macrophages:** Interact with COL11A1+ CAFs to form a physical barrier.
  - **Tumor Cells:** “Basal/Stress” state characterized by NRF2 pathway activation.
  - **T Cells:** Distinct exhaustion trajectories vs. activation.

### 2.2 Setup

```
library(clusterProfiler)
library(dplyr)
```

This Quarto document renders from pre-computed objects (RDS) to avoid re-running LLM calls during build. To reproduce the results end-to-end, you will need the development version of `clusterProfiler` that provides `interpret()`, `interpret_agent()`, and `interpret_hierarchical()`, and an LLM backend configured via `fanyi`.

```
library(fanyi)

api_key <- Sys.getenv("DEEPSEEK_API_KEY")
set_translate_option(source = "dsk", key = api_key)
```

## 3 Analysis Workflow

### 3.1 1. Data Preparation

We load the pre-processed marker genes and metadata.

The case study is distributed with the following inputs in the same directory as this QMD:

File	What it contains	Used for
<code>scobj.markers.rds</code>	FindAllMarkers output (per-cluster marker genes)	Selecting top markers for enrichment
<code>md.rds</code>	Cell-level metadata from the Seurat object	Building context and priors
<code>x.rds</code>	Cell-marker enrichment ( <code>compareCluster(..., fun=enricher)</code> )	LLM interpretation input
<code>anno.rds, pheno.rds, deep.rds, hier.rds</code>	Pre-computed LLM outputs	Rendering results without live API calls

The `md.rds` file is a cell-level metadata table (i.e., `scobj@meta.data`) saved separately so this case study does not need to bundle the full Seurat object. It is expected to contain, at minimum:

- Study/clinical fields used to summarize dataset context: `PathType`, `Timepoint`, `PathRes`, `Drug`, `Group`
- Cluster ID used to aggregate priors: `seurat_clusters`
- A coarse lineage/cell class label used as prior knowledge: `all_cluster_annotation`

```
scobj.markers <- readRDS("./scobj.markers.rds")

# Filter for top 10 markers per cluster
scobj.markers %>%
  group_by(cluster) %>%
  dplyr::filter(avg_log2FC > 1) %>%
  slice_head(n = 10) %>%
  ungroup() -> top10

md <- readRDS("./md.rds")
```

### 3.2 2. Context Construction

We construct a prompt context based on the study design.

```
majority <- function(x) names(sort(table(x), decreasing = TRUE))[1]

ctx <- paste0(
  "Study setting: NSCLC single-cell RNA-seq in the context of ",
  "neoadjuvant ICB + chemotherapy. ",
  "This dataset includes paired/longitudinal sampling ",
  "(pre- vs post-treatment) and response stratification by pathology. ",
  "Current object summary: ",
  "PathType=", majority(md$PathType), "; ",
  "Timepoint=", majority(md$Timepoint), "; ",
  "PathologicResponse=", majority(md$PathRes), "; ",
  "Response labels include NMPR/MPR/pCR. ",
```

```

"Drug=", majority(md$Drug), "; ",
"Group=", majority(md$Group), ". ",
"Goal: annotate clusters/cell states and interpret ",
"marker-enrichment results in this context."
)

```

We also construct two optional inputs used by the interpretation functions:

1. A prior vector (`prior_vec`) derived from the dataset's coarse lineage annotation (`all_cluster_annotation`).
2. A per-gene fold-change lookup (`gene_fold_change`) derived from marker statistics.

```

prior_tbl <- md %>%
  mutate(cluster = .data[["seurat_clusters"]]) %>%
  group_by(cluster) %>%
  summarise(prior = majority(all_cluster_annotation), .groups = "drop")
prior_vec <- setNames(prior_tbl$prior, as.character(prior_tbl$cluster))

fc_col <- if ("avg_log2FC" %in% names(top10)) "avg_log2FC" else "avg_logFC"
gene_fc <- top10 %>%
  group_by(gene) %>%
  summarise(fc = max(.data[[fc_col]], na.rm = TRUE), .groups = "drop")
gene_fold_change <- setNames(gene_fc$fc, gene_fc$gene)

```

### 3.3 3. Standard Enrichment Analysis

```
x <- readRDS("./x.rds")
```

The enrichment objects were generated from `top10` using `compareCluster`. Cell-marker enrichment used an external TERM2GENE table (`Cell_marker_Human.xlsx`), which is not bundled in this directory.

```

cm <- readxl::read_xlsx("./Cell_marker_Human.xlsx")

x <- compareCluster(
  gene ~ cluster,
  data = top10,
  fun = enricher,
  TERM2GENE = cm[, c("cell_name", "marker")]
)
saveRDS(x, "./x.rds")

```

## 4 Results: LLM-Driven Interpretation

### 4.1 Task 1: Cell Type Annotation (Hypothesis Verification)

We use `interpret(task = "celltype")` with prior knowledge injection.

```
library(fanyi)

api_key <- Sys.getenv("DEEPSEEK_API_KEY")
set_translate_option(source = "dsk", key = api_key)

anno <- interpret(
  x = x,
  context = ctx,
  task = "celltype",
  prior = prior_vec,
  n_pathways = 20,
  add_ppi = FALSE,
  gene_fold_change = gene_fold_change,
  model = "deepseek-chat",
  api_key = api_key
)
saveRDS(anno, "./anno.rds")
```

#### Result:

Full results (online): [https://yulab-smu.top/clusterProfiler-LLM-Supplementary/Task1\\_Cell-Type\\_Annotation.md](https://yulab-smu.top/clusterProfiler-LLM-Supplementary/Task1_Cell-Type_Annotation.md)

#### Comparison with *Nature Genetics* (2025):

Biological Feature	clusterProfiler-LLM Interpretation	Nature Genetics (Original Findings)	Key Match?	Significance
<b>Macrophage Identity</b>	Identified <b>SPI1+ M2-like Macrophages</b> with immunosuppressive signatures ( <i>C1QA/B/C, CD163</i> ).	Described <b>SPP1+ TAMs</b> as a major immunosuppressive population interacting with CAFs.	<b>Yes</b> (M2/TAM overlap)	Major Resistance Mechanism

Biological Feature	clusterProfiler-LLM Interpretation	Nature Genetics (Original Findings)	Key Match?	Significance
<b>Tumor Cell State</b>	Detected <b>Basal-like Progenitor</b> state with <i>GPX2</i> and <i>AKR1C1</i> (Oxidative Stress).	Defined “ <b>Basal/Stress</b> ” tumor cell state characterized by <b>NRF2</b> activation.	<b>Direct</b> ( <b>Stress/Basal</b> )	Poor Prognosis Driver
<b>Stromal Niche</b>	Annotated <b>Activated Myofibroblasts</b> ( <i>ACTA2</i> , <i>COL1A2</i> ) and Vascular Progenitors.	Detailed <b>COL1A1+</b> CAFs forming a barrier at tumor boundaries.	<b>Yes</b> (CAF Activation)	Immune Exclusion
<b>T Cell States</b>	Found <b>IL7R+ Memory T cells</b> and <b>Exhausted</b> subsets.	Highlighted dysfunctional CD8+ T cells and potential for revival.	<b>Yes</b>	Therapy Response Determinant

## 4.2 Task 2: Phenotypic Characterization

We use `interpret(task = "phenotype")` to define the functional state of each cluster.

```
library(fanyi)

api_key <- Sys.getenv("DEEPSEEK_API_KEY")
set_translate_option(source = "dsk", key = api_key)

pheno <- interpret(
  x = x,
  context = ctx,
  task = "phenotype",
  n_pathways = 30,
  add_ppi = FALSE,
  gene_fold_change = gene_fold_change,
  model = "deepseek-chat",
  api_key = api_key
)
saveRDS(pheno, "./pheno.rds")
```

**Result:**

Full results (online): [https://yulab-smu.top/clusterProfiler-LLM-Supplementary/Task2\\_Phenotypic\\_Characterization.md](https://yulab-smu.top/clusterProfiler-LLM-Supplementary/Task2_Phenotypic_Characterization.md)

### Comparison with *Nature Genetics* (2025):

Biological Feature	clusterProfiler-LLM Interpretation	Nature Genetics (Original Findings)	Key Match?	Significance
<b>Tumor Cell Plasticity</b>	Identified <b>Basal-like Progenitor</b> state expressing <i>AKR1C1, GPX2</i> (Detoxification/Stress).	Defined “ <b>Basal/Stress</b> ” state driven by <i>NRF2 (NFE2L2)</i> and oxidative stress response.	<b>Direct</b> (AKR1C1/Stress)	Chemosensitivity Mechanism
<b>Fibroblast Activation</b>	Characterized <b>Activated Myofibroblasts</b> with high <i>COL1A1, ACTA2, TGFB1</i> .	Described <b>COL1A1+</b> CAFs forming immune-excluded niches.	<b>Yes</b> (TGF-( ) CAFs)	Immune Exclusion Barrier
<b>Myeloid Polarization</b>	Detected <b>Immunosuppressive Macrophages</b> ( <i>C1QA, APOE</i> ) with M2-like features.	Highlighted <b>SPP1+ TAMs</b> as the dominant immunosuppressive myeloid population.	<b>Yes</b> (M2/TAMs)	T-cell Suppression
<b>B Cell Function</b>	Noted <b>Germinal Center B cells</b> ( <i>BCL6, AICDA</i> ) indicating humoral immunity.	Discussed <b>TLS (Tertiary Lymphoid Structure)</b> presence and B cell maturity.	<b>Consistent</b>	Anti-tumor Immunity Potential

### 4.3 Task 3: Mechanism Interpretation (Multi-Agent Deep Mode)

We use `interpret_agent()` to reconstruct causal networks, integrating PPI and LogFC data.

```
library(fanyi)

api_key <- Sys.getenv("DEEPSEEK_API_KEY")
set_translate_option(source = "dsk", key = api_key)

deep <- interpret_agent(
  x = x,
  context = ctx,
```

```

n_pathways = 50,
add_ppi = TRUE,
gene_fold_change = gene_fold_change,
model = "deepseek-chat",
api_key = api_key
)
saveRDS(deep, "./deep.rds")

```

### Result:

Full results (online): [https://yulab-smu.top/clusterProfiler-LLM-Supplementary/Task3\\_Mechanism\\_Interpretation.md](https://yulab-smu.top/clusterProfiler-LLM-Supplementary/Task3_Mechanism_Interpretation.md)

### Comparison with *Nature Genetics* (2025):

Biological Feature	clusterProfiler-LLM Interpretation	Nature Genetics (Original Findings)	Key Match?	Significance
<b>T Cell Exhaustion</b>	<b>TOX</b> identified as a key regulator of <b>Exhausted CD8+ T cells</b> ( <i>PDCD1</i> , <i>HAVCR2</i> ).	<b>Dysfunctional CD8+ T cells</b> are a major feature of non-responders; exhaustion is a key barrier.	<b>Direct</b> ( <b>TOX/Exhaustion</b> )	Checkpoint Blockade Target
<b>Myeloid Regulation</b>	<b>SPI1 (PU.1)</b> inferred as the master regulator for <b>M2-like TAMs</b> ( <i>C1QA</i> , <i>MRC1</i> ).	<b>SPP1+ TAMs</b> recruit regulatory T cells and suppress adaptive immunity.	<b>Yes</b> (Myeloid Driver)	Immunosuppressive Hub
<b>B Cell Identity</b>	<b>PAX5</b> and <b>BCL6</b> regulatory network defines <b>Follicular/GC B cells</b> .	Presence of mature B cells in <b>TLS</b> correlates with better prognosis.	<b>Consistent</b>	Prognostic Marker
<b>Tumor Proliferation</b>	<b>E2F</b> targets and <b>MYC</b> signaling active in <b>Cycling Tumor Cells</b> .	High proliferation rates in specific tumor subclones drive progression.	<b>Yes</b>	Tumor Aggressiveness

### 4.4 Task 4: Hierarchical Interpretation

Refining annotations from Major lineage to Minor states using `interpret_hierarchical()`.

```

library(fanyi)

api_key <- Sys.getenv("DEEPSEEK_API_KEY")
set_translate_option(source = "dsk", key = api_key)

map_tbl <- md %>%
  mutate(minor = as.character(minor), major = as.character(major)) %>%
  count(minor, major, name = "n") %>%
  group_by(minor) %>%
  slice_max(n, n = 1, with_ties = FALSE) %>%
  ungroup()
mapping <- setNames(map_tbl$major, map_tbl$minor)

top10_minor <- readRDS("./top10_minor.rds")
top10_major <- readRDS("./top10_major.rds")

cm <- readxl::read_xlsx("./Cell_marker_Human.xlsx")

x_minor <- compareCluster(
  gene ~ cluster,
  data = top10_minor,
  fun = enricher,
  TERM2GENE = cm[, c("cell_name", "marker")]
)
saveRDS(x_minor, "./x_minor.rds")

x_major <- compareCluster(
  gene ~ cluster,
  data = top10_major,
  fun = enricher,
  TERM2GENE = cm[, c("cell_name", "marker")]
)
saveRDS(x_major, "./x_major.rds")

hier <- interpret_hierarchical(
  x_minor = x_minor,
  x_major = x_major,
  mapping = mapping,
  model = "deepseek-chat",
  api_key = api_key,
  task = "cell_type"
)
saveRDS(hier, "./hier.rds")

```

## Result:

Full results (online): [https://yulab-smu.top/clusterProfiler-LLM-Supplementary/Task4\\_Hierarchical\\_Interpretation.md](https://yulab-smu.top/clusterProfiler-LLM-Supplementary/Task4_Hierarchical_Interpretation.md)

### Comparison with *Nature Genetics* (2025):

Biological Feature	clusterProfiler-LLM Interpretation	Nature Genetics (Original Findings)	Key Match?	Significance
<b>CD8+ T Cell Subsets</b>	Distinctly separated <b>Naïve</b> , <b>Memory</b> ( <i>IL7R</i> ), and <b>Exhausted</b> ( <i>TOX</i> , <i>LAG3</i> ) T cells.	Emphasized the spectrum from <b>Pre-dysfunctional</b> to <b>Dysfunctional</b> CD8+ T cells.	<b>Precise</b>	Therapy Response Continuum
<b>Macrophage Heterogeneity</b>	Hierarchically resolved <b>Alveolar Macrophages</b> vs. <b>Tumor-Associated Macrophages</b> ( <i>C1QA</i> ).	Distinguished resident macrophages from tumor-infiltrating <b>SPP1+</b> TAMs.	<b>Yes</b>	Origin Matters (Tissue vs. Tumor)
<b>Tumor Heterogeneity</b>	Sub-classified tumor cells into <b>Cycling</b> ( <i>MKI67</i> ) and <b>Stress/Basal</b> ( <i>AKR1C1</i> ) states.	Mapped tumor macro-clusters to distinct <b>Cellular States</b> (Cycling, Stress, Interferon-high).	<b>Direct</b>	Intratumoral Heterogeneity
<b>Endothelial States</b>	Identified <b>Tip cells</b> and <b>Stalk cells</b> indicative of angiogenesis.	Noted <b>PLVAP+</b> <b>Endothelial cells</b> associated with tumor vascularization.	<b>Consistent</b>	Angiogenesis Targets

## 5 Conclusion

The clusterProfiler-LLM framework successfully reproduced the key biological insights of the 2025 *Nature Genetics* study without manual curation. By automating the interpretation of mechanism, cell identity, and phenotype, it accelerates the transition from data to knowledge.