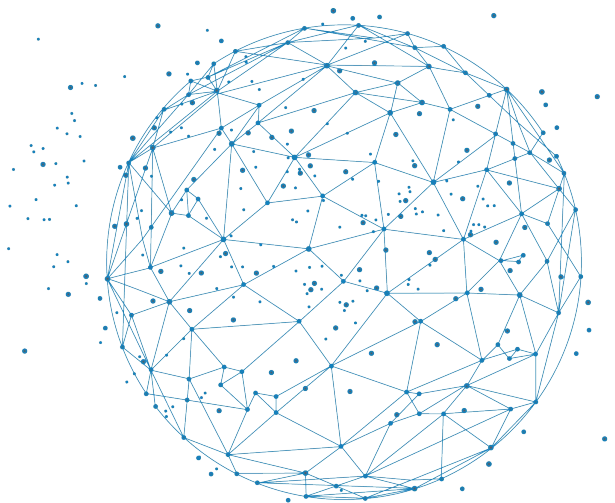




Third edition 2026 in Fréjus

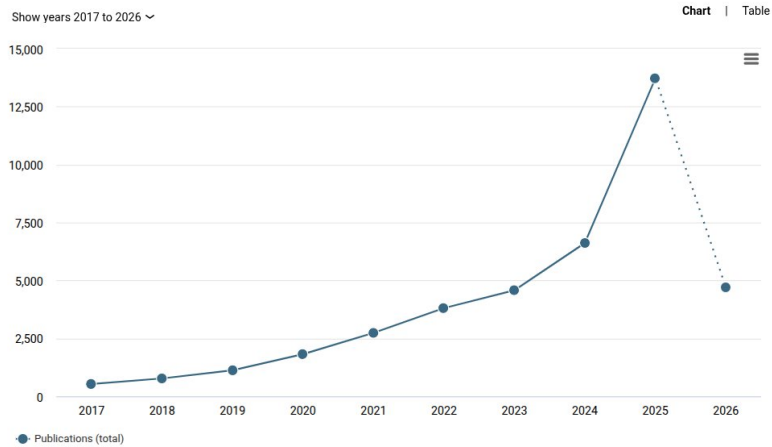


Omics integration - General aspects

Jimmy Vandel



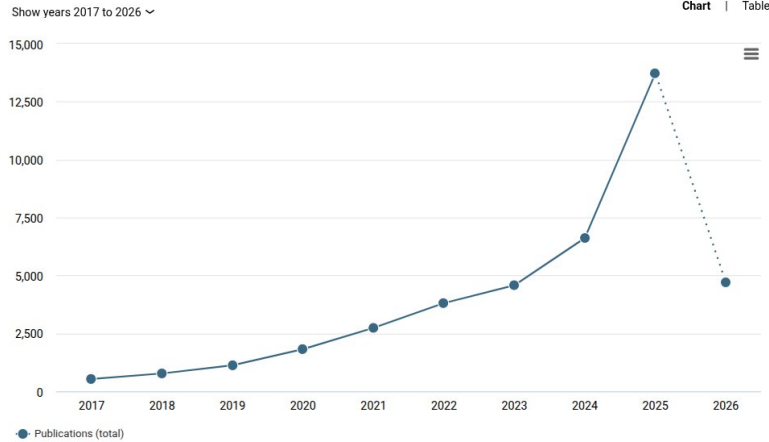
“Multi-omics” citations



<https://app.dimensions.ai/discover/publication> (22th Mar. 2026: 163,356,634 referenced publications)

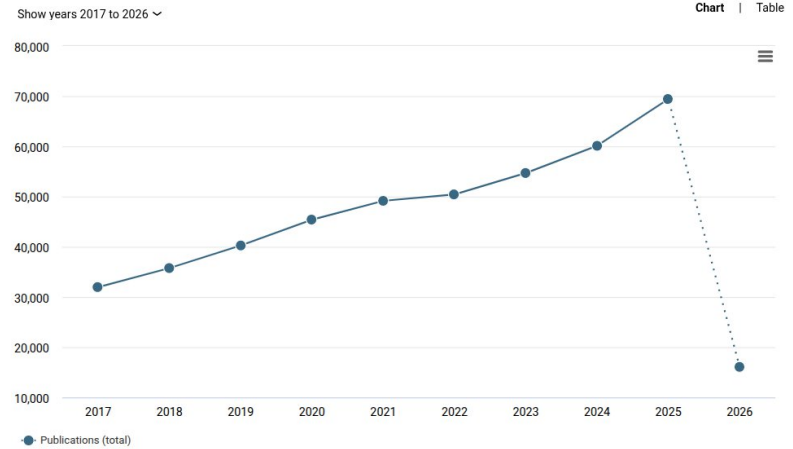


“Multi-omics” citations



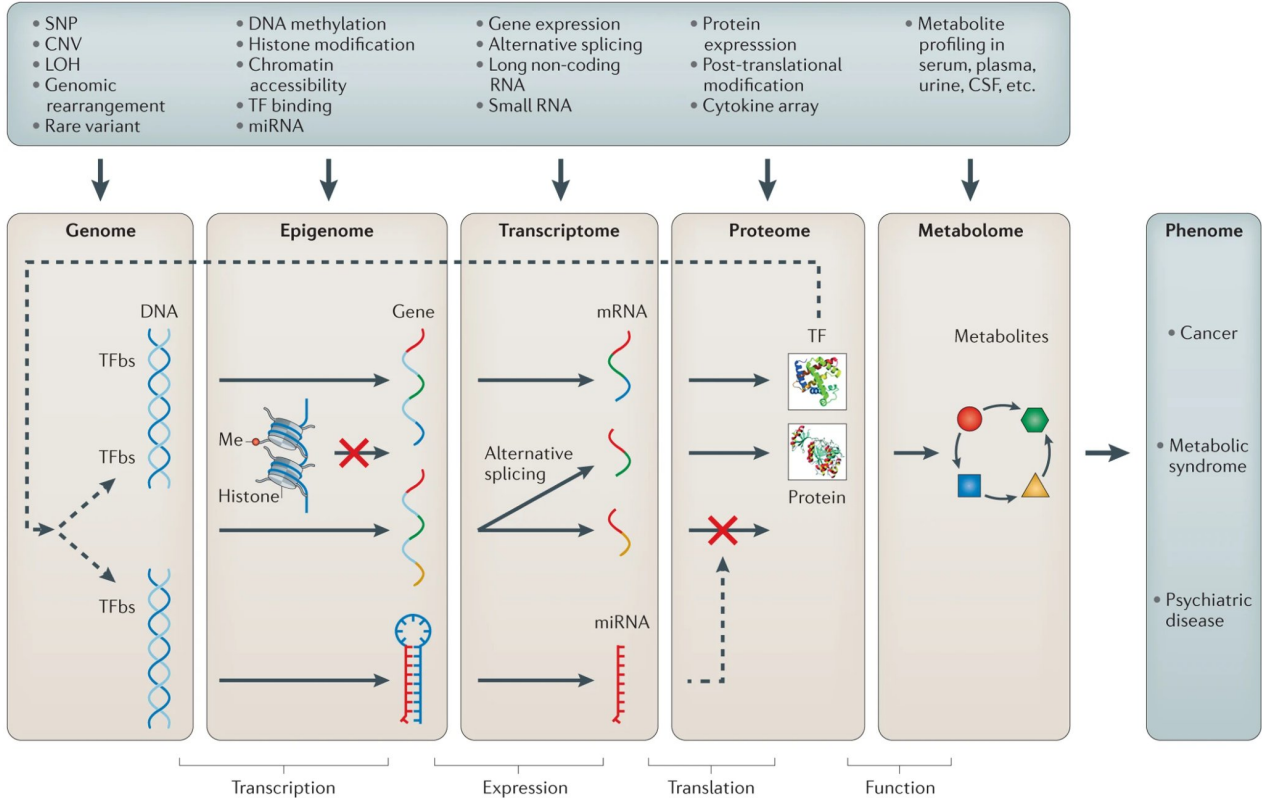
| Name | ↓ Publications | Citations |
|--|----------------|-----------|
| Fields of Research (ANZSRC 2020) code | | |
| Biomedical and Clinical Sciences 32 | 22,930 | 263,806 |
| Biological Sciences 31 | 20,195 | 348,609 |
| Oncology and Carcinogenesis 3211 | 7,755 | 73,484 |
| Medical Biochemistry and Metabolomics 3205 | 6,286 | 90,509 |
| Genetics 3105 | 6,141 | 111,426 |
| Bioinformatics and Computational Biology 3102 | 5,764 | 105,554 |

“Single-cell” citations



| Name | ↓ Publications | Citations |
|--|----------------|------------|
| Fields of Research (ANZSRC 2020) code | | |
| Biomedical and Clinical Sciences 32 | 515,971 | 18,438,401 |
| Biological Sciences 31 | 329,178 | 14,987,432 |
| Oncology and Carcinogenesis 3211 | 164,331 | 4,548,786 |
| Biochemistry and Cell Biology 3101 | 164,230 | 7,806,708 |
| Engineering 40 | 141,413 | 4,130,923 |
| Clinical Sciences 3202 | 121,843 | 3,710,657 |

Omics... which ones ?



Nature Reviews | **Genetics**

Ritchie, M., Holzinger, E., Li, R. et al. *Methods of integrating data to uncover genotype-phenotype interactions.* Nat Rev Genet 16, 85-97 (2015).

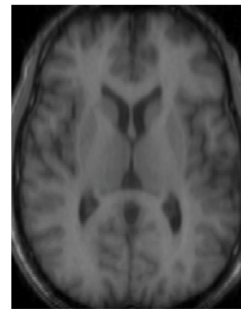


Other related data ?

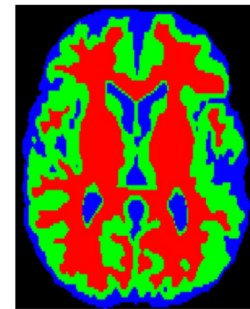
- clinical data
- imaging data / radiomics / pathomics (full data or extracted characteristics)
- new omics fields : fluxomics, ionomics, microbiomics, glycomics...
- biological knowledge : DNA/protein, protein/protein interactions, DNA recombination
→ a priori in model definition/construction



| CBC Information | | | | | | | | | | Other Tests: | | | | Click on a ? for additional help information | | |
|-----------------|------|------|------|-----|-----------|----------------|-----------------|---------------|----------------|--------------|---------|--------|-------|---|---------|----------|
| Date | WBC | RBC | HGB | HCT | Platelets | Percent Lymphs | Absolute Lymphs | Percent Neuts | Absolute Neuts | B/M | BCV | BCM | BCDC | | RDW | SPV |
| ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | | ? | ? |
| 20-Jan-15 | 9.0 | 5.00 | 9.0 | 45 | 400 | 70.0% | 3.5 | 45.0% | 2.3 | 48.0% | 2.2 | 80.100 | 29.32 | 32.35 | 11% 10% | 7.5 11 5 |
| 20-Jan-15 | 19.0 | 4.80 | 9.5 | 42 | 400 | 50.0% | 5.8 | 55.0% | 5.5 | 5.0 | 1.2 2.0 | | | | | |
| 20-Jan-16 | 12.0 | 5.10 | 10.0 | 38 | 400 | 75.0% | 7.7 | 50.0% | 6.0 | 2.0 | | | | | | |
| 20-Jan-16 | 11.0 | 5.20 | 12.0 | 33 | 400 | 68.0% | 6.9 | 45.0% | 5.0 | 2.0 | | | | | | |
| 20-Jan-17 | 8.0 | 5.00 | 11.0 | 34 | 400 | 68.0% | 5.2 | 40.0% | 3.2 | 2.0 | | | | | | |
| 20-Jan-17 | 7.0 | 5.30 | 13.0 | 32 | 400 | 70.0% | 4.2 | 40.0% | 2.8 | 2.0 | | | | | | |
| 20-Jan-18 | 5.0 | 5.40 | 15.0 | 30 | 400 | 70.0% | 3.5 | 45.0% | 2.3 | 2.0 | | | | | | |
| 20-Jan-18 | 4.5 | 5.80 | 13.8 | 40 | 250 | 50.0% | 2.3 | 48.0% | 2.2 | 2.0 | | | | | | |
| 20-Jan-19 | 4.0 | 6.00 | 14.0 | 48 | 150 | 75.0% | 3.0 | 50.0% | 2.0 | 2.0 | | | | | | |
| 20-Jan-19 | 7.0 | 5.00 | 12.0 | 45 | 140 | 65.0% | 4.6 | 51.0% | 3.6 | 2.0 | | | | | | |
| 20-Jan-20 | 9.0 | 4.50 | 10.0 | 47 | 130 | 68.0% | 6.1 | 55.0% | 5.0 | 2.0 | | | | | | |
| 20-Jan-20 | 10.0 | 5.20 | 11.0 | 45 | 250 | 55.0% | 5.5 | 60.0% | 6.0 | 2.0 | | | | | | |



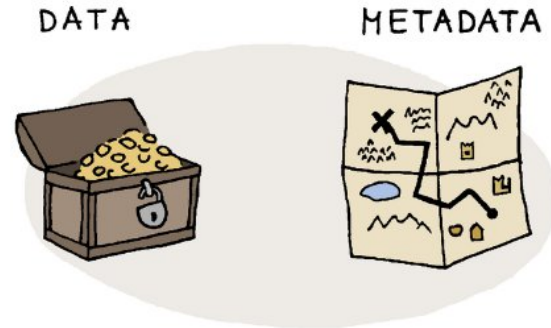
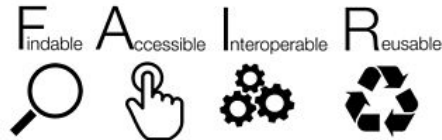
(a) Axial slice



(b) Tissue segmentation

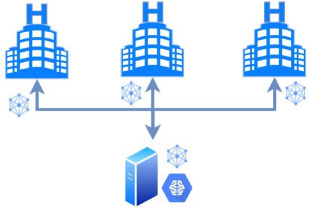


- **Omics data:** mostly dependent on omics and/or technology (FASTQ, VCF, mzML...)
→ numeric matrices
- **Clinical:** very heterogeneous (care/research), coding standards (diagnoses, medication...), electronic health records with interoperable standards such as openEHR or FHIR
→ numeric matrices
- **Imaging data:** Digital Imaging and Communications in Medicine (DICOM) standard (image storage, metadata definition and communication protocols) but also TIFF, NifTI
→ numeric matrices



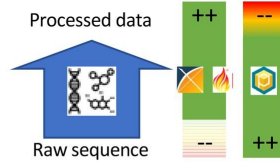


GENOMED4ALL

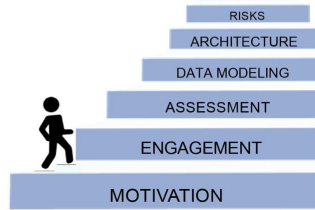


Challenges

Assess interoperability standards



Effective data modeling strategy

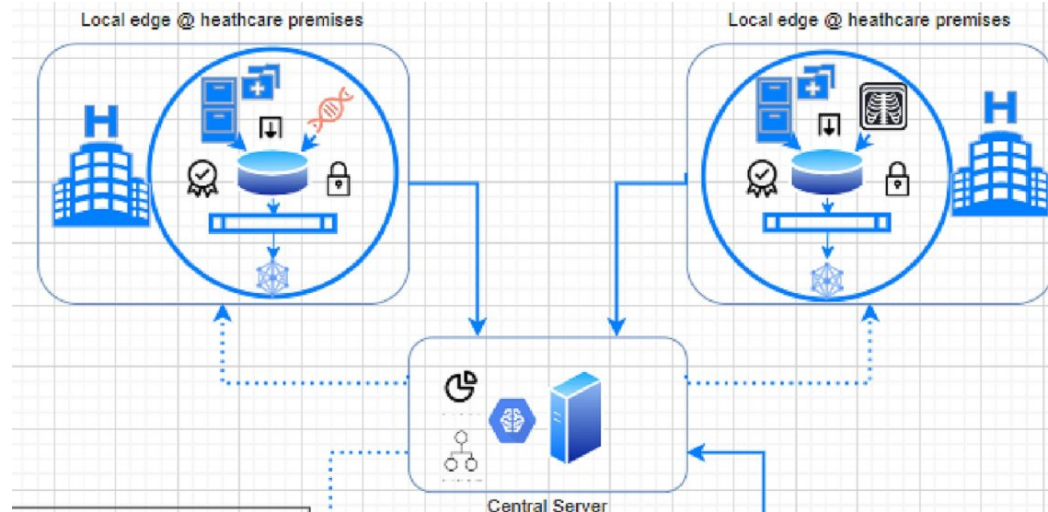


Common data model for federated learning on multi-omics and clinical data?



Federated learning

Cremonesi F et al.. The need for multimodal health data modeling: A practical approach for a federated-learning healthcare platform. *J Biomed Inform.* 2023



Génomique

- WGS
- Génotypage (puces ADN, PMRA)
- Variants génétiques (eQTL, mQTL, GWAS)
- Métagénomique (16S, short/long reads)

Transcriptomique

- RNA-seq (bulk, long/short reads, Nanopore, Illumina)
- scRNA-seq/snRNA-seq (10X Genomics, Visium, etc.)
- Nascent RNA-seq
- Métatranscriptomique (microbiote, pathogènes)
- miRNA-seq (puces ou séquençage)
- CITE-seq

Épigénomique

- Méthylation (EPIC, bisulfite sequencing, RRBS, EM-seq)
- Accessibilité de la chromatine (ATAC-seq)
- Modifications d'histones (ChIP-seq, Cut&Run, Cut&Tag)
- Épitranscriptomique (m6A, R-loops)

Protéomique, métabolomique et lipidomique

- Spectrométrie de masse
- Ribosome profiling (traductome)
- Métaprotéomique (microbiote)
- Protéines associées à la chromatine (Pol II, etc.)

Mais aussi : données spatiales, imagerie, cytométrie, phénotypique, clinique

Analyse single-omics :

DESeq2, edgeR, Seurat, Boruta (métabolomique),
ChromHMM

Analyse d'enrichissement :

GO terms

Outils d'intégration :

WNN, WGCNA, SNF, iCluster, PLS-DA, MixOmics,
DIABLO, RGCCA, SGCCA, ComDim/CCSWA,
MOFA, network diffusion hierarchique, MINI-EX

Corrélation et clustering

Visualisation :

ShinyR, PCA/MDS (réduction de
dimension)

Analyse single-omics :

DESeq2, edgeR, Seurat, Boruta (métabolomique),
ChromHMM

Analyse d'enrichissement :

GO terms

Outils d'intégration :

WNN, WGCNA, SNF, iCluster, PLS-DA, MixOmics,
DIABLO, RGCCA, SGCCA, ComDim/CCSWA,
MOFA, network diffusion hierarchique, MINI-EX

Corrélation et clustering

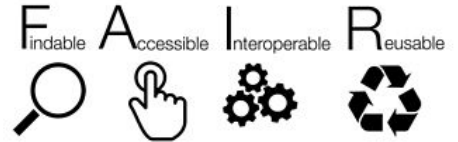
Visualisation :

ShinyR, PCA/MDS (réduction de
dimension)





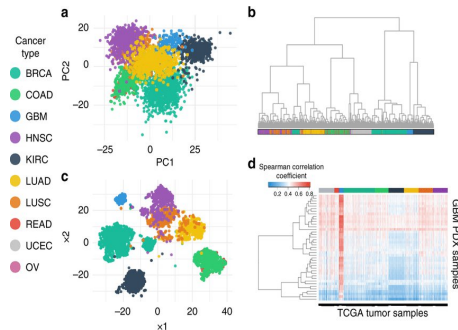
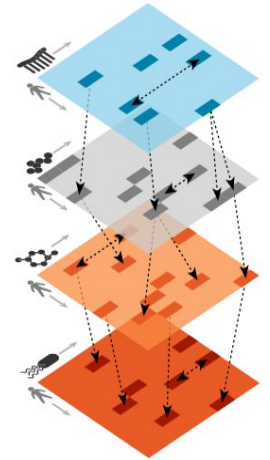
- Take advantage of the vast amount of available data
 - Data access (local/national regulation, infrastructures...)
 - Data representation (structuration, ontologies...)→ Need of common representation framework



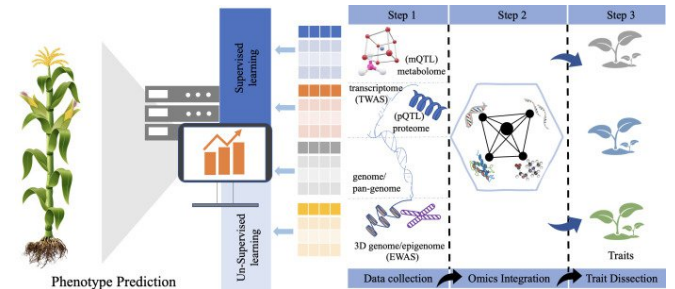
- Improve our understanding of biological phenomena
 - Data heterogeneity (technology, format, meaning, distribution...)
 - Data complexity (dependances/independances, ad-hoc assumptions...)
 - Amount a data (time/memory consuming)→ Need for methods adapted to these data



- Deep insights into biology phenomenon
- Subtyping and classification (disease, species, varieties)
- Biomarkers prediction: diagnostic, disease drivers, plant/animal selection...



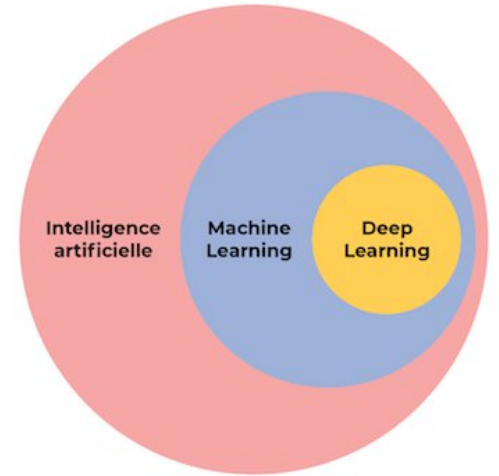
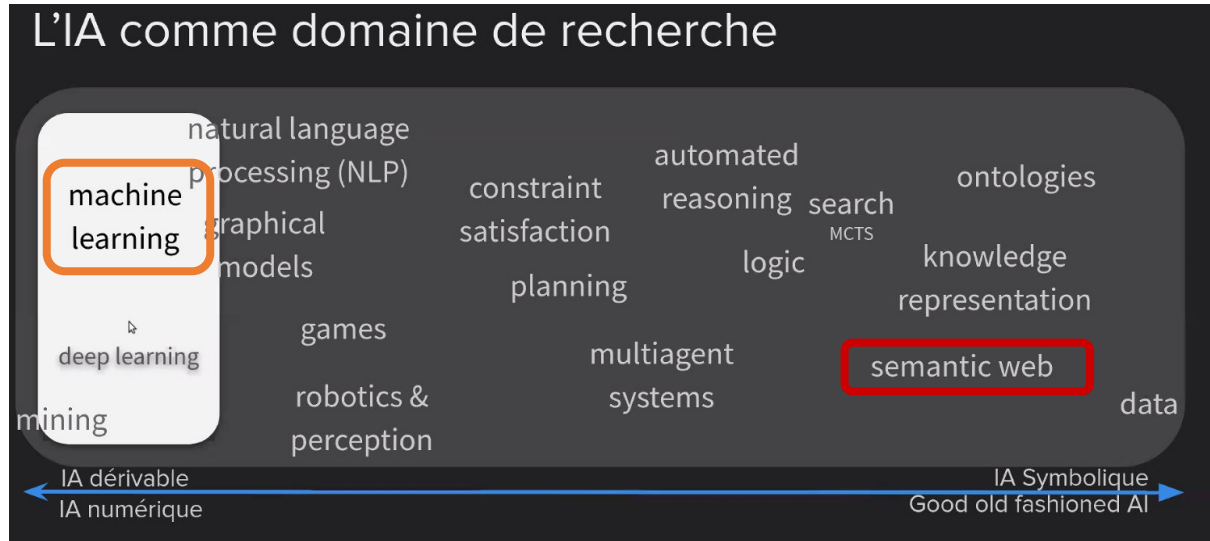
Vasileios et al (2018). Drug and disease signature integration identifies synergistic combinations in glioblastoma. Nature Communications. 9.



Mahmood et al (2022) Multi-omics revolution to promote plant breeding efficiency. Front Plant Sci Dec 8.



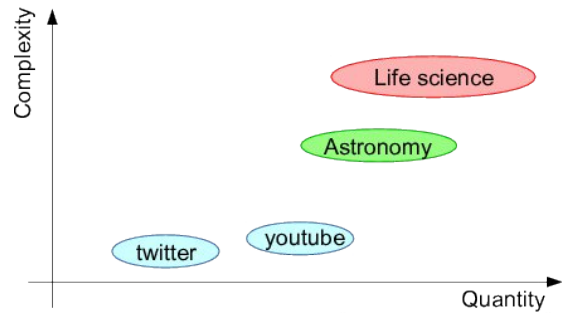
Artificial intelligence of course ... and so ?



Improve our understanding of biological phenomena

Take advantage of the vast amount of available data

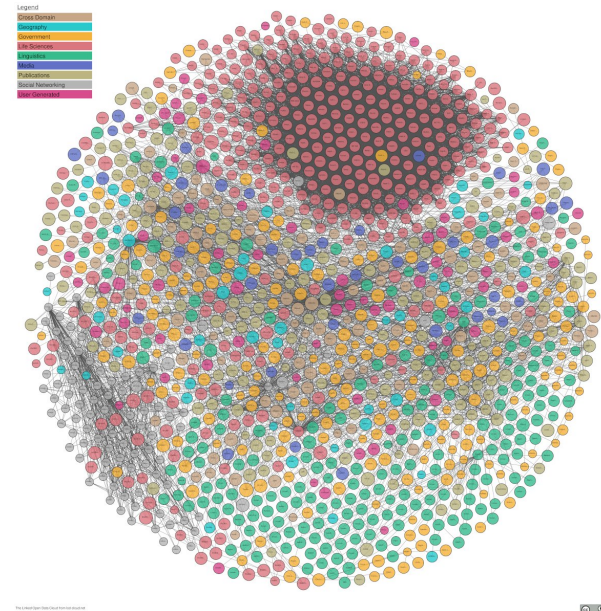
Take advantage of the vast amount of available data



Big Data: Astronomical or Genomical?

Zachary D. Stephens¹, Skylar Y. Lee¹, Faraz Faghri², Roy H. Campbell², Chengxiang Zhai³, Miles J. Efron⁴, Ravishankar Iyer¹, Michael C. Schatz^{5*}, Saurabh Sinha^{3*}, Gene E. Robinson^{6*}

PLOS Biology | DOI:10.1371/journal.pbio.1002195 July 7, 2015



Life science: 1600+ reference databases

→ integrating heterogeneous data and knowledge is (badly) needed!

Editorial > Nucleic Acids Res. 2022 Jan 7;50(D1):D1-D10. doi: 10.1093/nar/gkab1195.

The 2022 Nucleic Acids Research database issue and the online molecular biology database collection

Daniel J Rigden¹, Xosé M Fernández²

Affiliations + expand

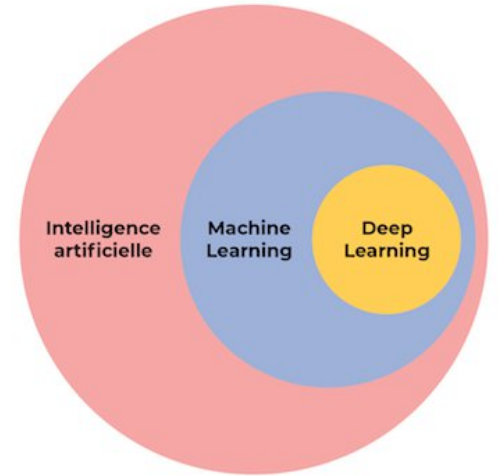
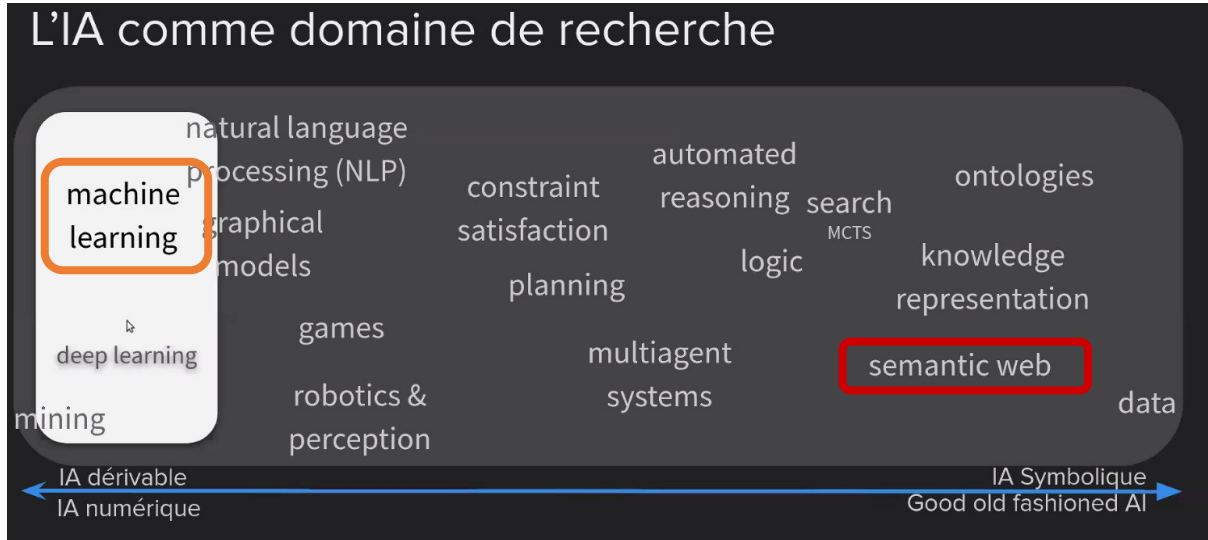
PMID: 34986604 PMCID: PMC8728296 DOI: 10.1093/nar/gkab1195

Semantic Web = framework for:

- integrating data and knowledge
- querying
- reasoning

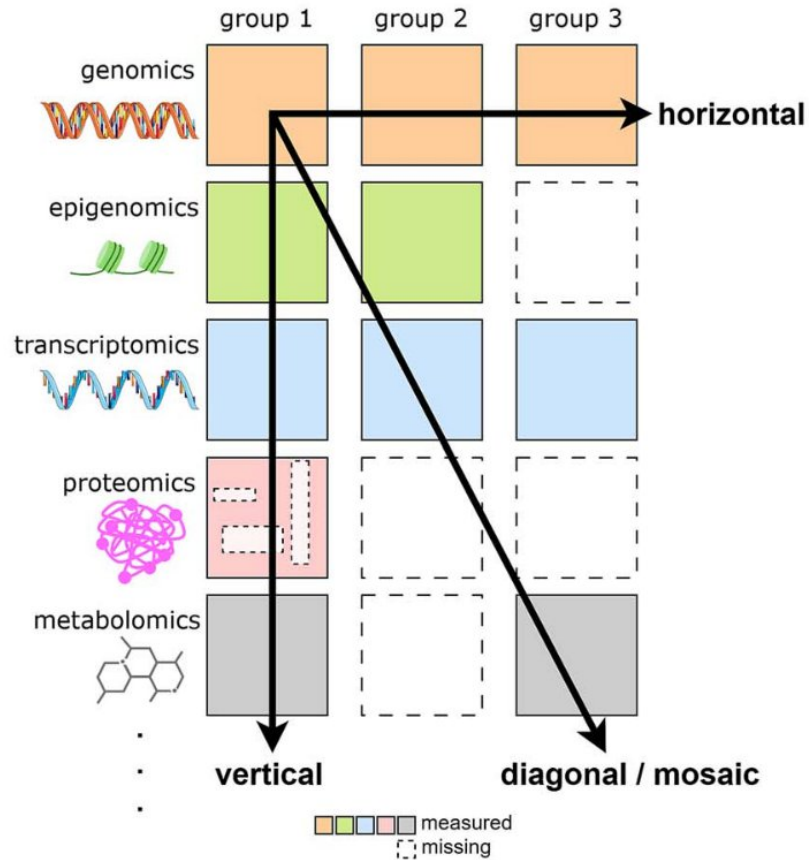


Artificial intelligence of course ... and so ?



Improve our understanding of biological phenomena

Take advantage of the vast amount of available data



- **Vertical integration:** combines different omics modalities within the same group of samples
- **Horizontal integration:** aligns datasets from the same omics layer across different sample groups (e.g. batches, cellular models)
- **Diagonal integration:** combines distinct omics modalities from different sample groups to explore inter-modality relationships across groups
- **Mosaic integration:** leverages overlapping modalities across samples to infer relationships and impute missing modalities



– Unsupervised learning

find hidden patterns, analyze and organize unlabelled data.

eg: clustering, dimension reduction, density estimation

- Self-supervised learning : optimise a loss function based on a ground truth

– Supervised learning

use labeled data to predict new observation outcome (predictive model).

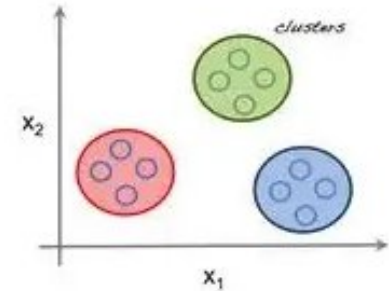
eg: classification task (categorical/numerical), regression (numerical)

– Semi-supervised

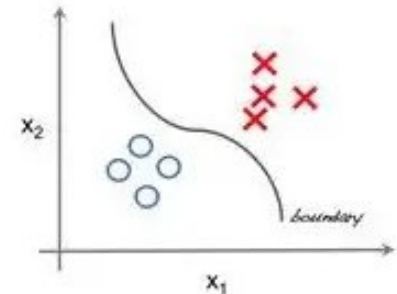
use labeled and unlabelled data to predict the unlabelled data and/or outcome

eg: inductive/transductive approaches

Unsupervised learning



Supervised learning

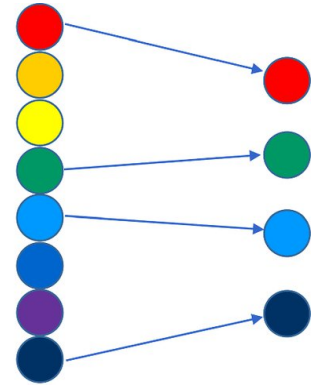




Feature selection

→ determine a smaller set of features minimizing (relevant) information loss

eg: filtering methods (correlation), recursive elimination, regularization

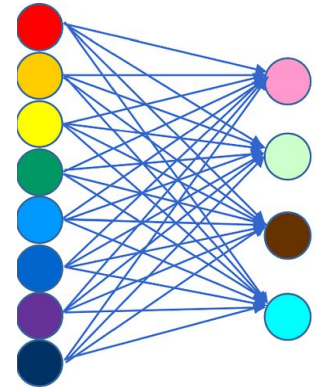


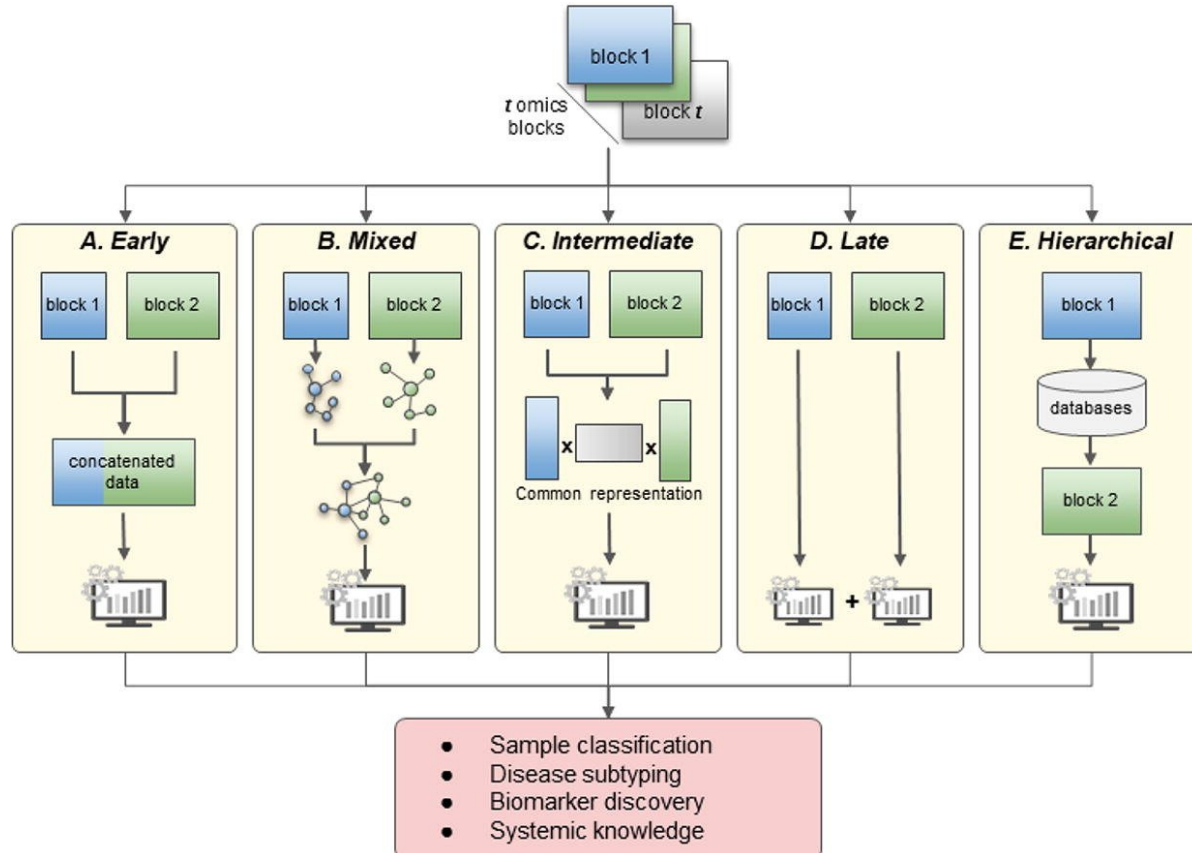
Feature extraction

→ combine the input features into another set of variables in a linear or non-linear fashion

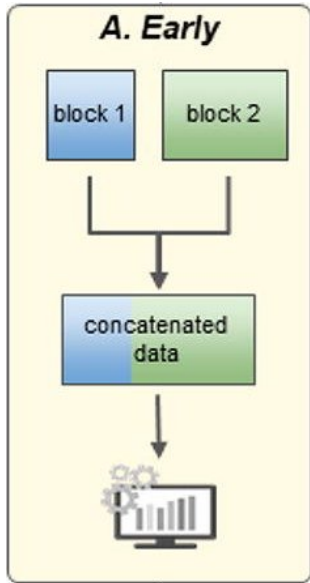
eg: PCA, PCoA, ICA...

+ regularization for sparse methods: sPCA, sNMF





Picard M. et al. Integration strategies of multi-omics data for machine learning analysis. *Comput Struct Biotechnol J.* 2021.



Concatenate every omics datasets into a single large matrix.

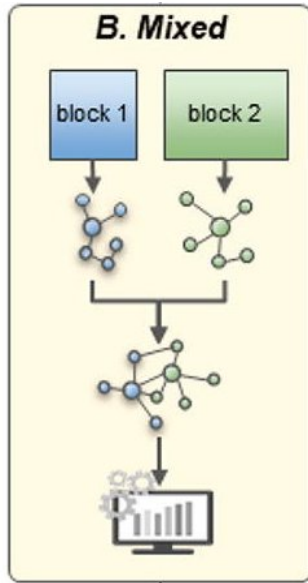
Pros :

- conceptually simple
- easy to implement
- can reveal interactions between omics

Cons :

- increased noise and dimensionality (concatenated matrix)
- need to deal with imbalanced omics datasets
- ignores the specific data distribution of each omics
- need a common dimension (rows or columns → samples or features)

Picard M. et al. Integration strategies of multi-omics data for machine learning analysis. *Comput Struct Biotechnol J.* 2021.



Transform independently each omics dataset into a simpler representation before integration.

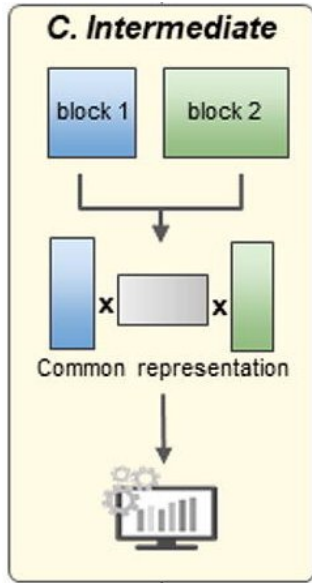
Pros :

- the new representation is less dimensional and less noisy
- reduce heterogeneity between omics
- classical approaches can be used on the new representation

Cons :

- the choice of the transformation method is not trivial
- risk of information loss during transformation
- need a common dimension (rows or columns → samples or features)

Picard M. et al. Integration strategies of multi-omics data for machine learning analysis. *Comput Struct Biotechnol J.* 2021.



Jointly integrate the multi-omics datasets without prior transformation.

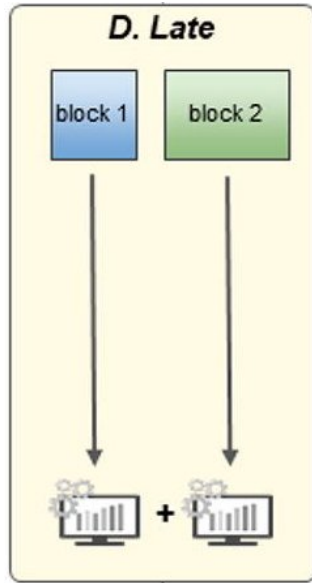
Pros :

- reduce information loss compare to the mixed strategy
- discover a joint inter-omics structure
- highlight the complementary information between omics

Cons :

- can require robust pre-processing step to reduce heterogeneity
- common latent space assumption

Picard M. et al. Integration strategies of multi-omics data for machine learning analysis. *Comput Struct Biotechnol J.* 2021.



Apply machine learning models separately on each omics dataset and then combine results.

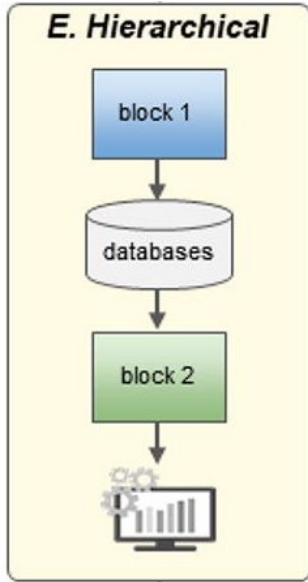
Pros :

- avoid the (many) challenges associated with direct omics integration
- you can use tools designed specifically for each omics
- classical approaches can be used to combine results

Cons :

- cannot capture direct inter-omics interactions
- complementarity information between omics is not exploited

Picard M. et al. Integration strategies of multi-omics data for machine learning analysis. *Comput Struct Biotechnol J.* 2021.



Include prior knowledge of omics relationships.

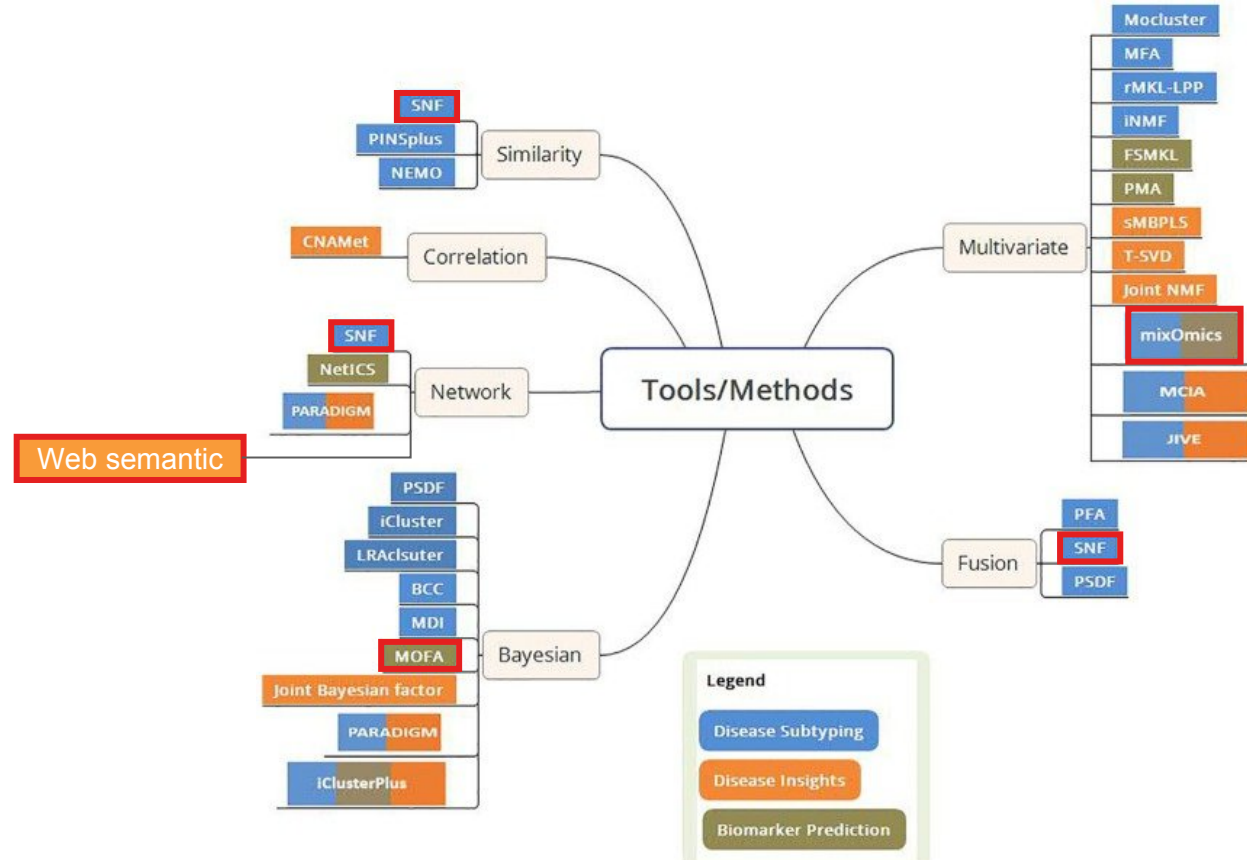
Pros :

- reduced complexity (sequential integration)
- integrate external knowledge

Cons :

- less generic than previous strategies

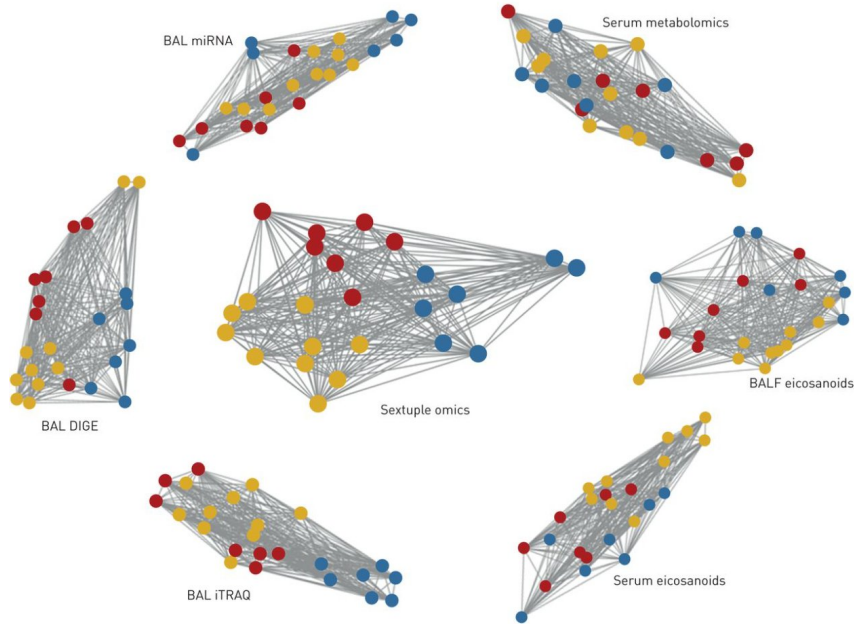
Picard M. et al. Integration strategies of multi-omics data for machine learning analysis. *Comput Struct Biotechnol J.* 2021.



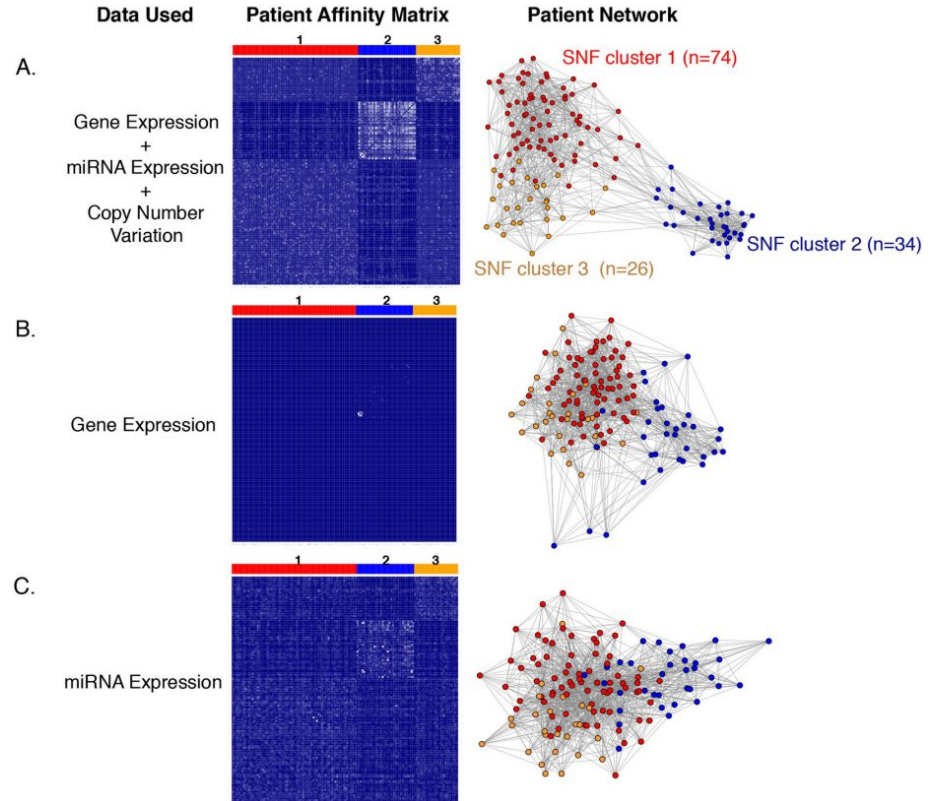
Subramanian I, Verma S, Kumar S, Jere A, Anamika K. Multi-omics Data Integration, Interpretation, and Its Application. *Bioinform Biol Insights*. 2020



Capturing shared topological structures and interactions across omics layers through networks

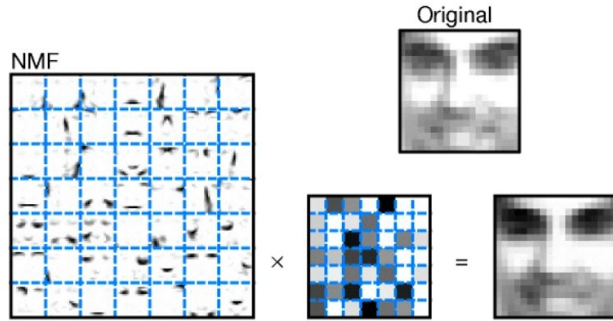


Li C et al. Integration of multi-omics datasets enables molecular classification of COPD. *Eur Respir J.* 2018

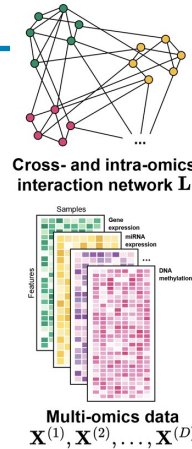
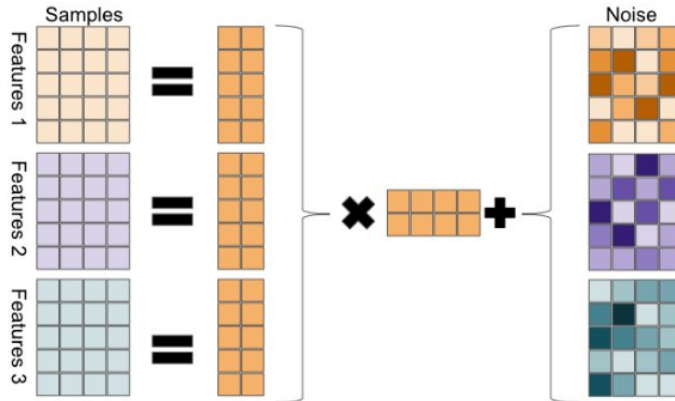


Chiu AM et al. Integrative analysis of the inter-tumoral heterogeneity of triple-negative breast cancer. *Sci Rep.* 2018

Matrix factorization based methods

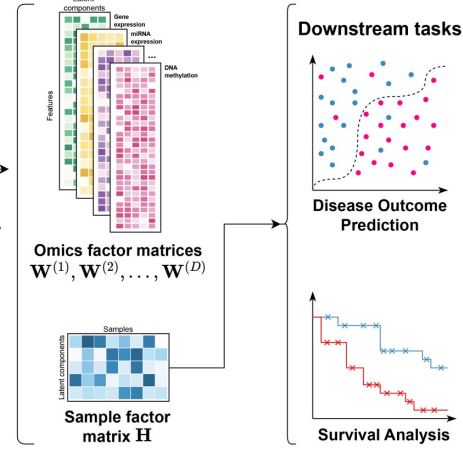


Lee, D et al. Learning the parts of objects by non-negative matrix factorization. Nature 1999.



X-intNMF

$$\min_{\mathbf{H}, \mathbf{W}^{(d)}} f = \frac{1}{2} \sum_{d=1}^D \|\mathbf{X}^{(d)} - \mathbf{W}^{(d)} \mathbf{H}\|_F^2 + \frac{\alpha}{2} \sum_{p=1}^D \sum_{q=1}^D \text{Tr}(\mathbf{W}^{(p)T} \mathbf{L}^{(pq)} \mathbf{W}^{(q)}) + \sum_{d=1}^D \beta_d \|\mathbf{W}^{(d)}\|_1 + \sum_{i=1}^N \gamma_i \|\mathbf{H}_{:,i}\|_1$$



Tien-Thanh Bui, et al. X-intNMF: a cross- and intra-omics regularized NMF framework for multi-omics integration, Bioinformatics 2026

<https://compgenomr.github.io/book/matrix-factorization-methods-for-unsupervised-multi-omics-data-integration.html>



Like PCA, a projection from the initial space onto a low-dimensional k subspace

$$y_{ijt} = \mathbf{x}_i \mathbf{\Gamma}_{jt} \boldsymbol{\beta}_{jt} + \varepsilon_{ijt}, \quad i = 1, \dots, n, j = 1, \dots, p_t, t \in (1, \dots, m)$$

$\mathbf{x}_i = (1, \mathbf{z}_i) = (1, z_{i1}, \dots, z_{ik})$ *projection vector in the k -subspace*

$\mathbf{\Gamma}_{jt} = \text{diag}(1, \gamma_{jt}, \dots, \gamma_{jt})$ *indicative diagonal matrix*

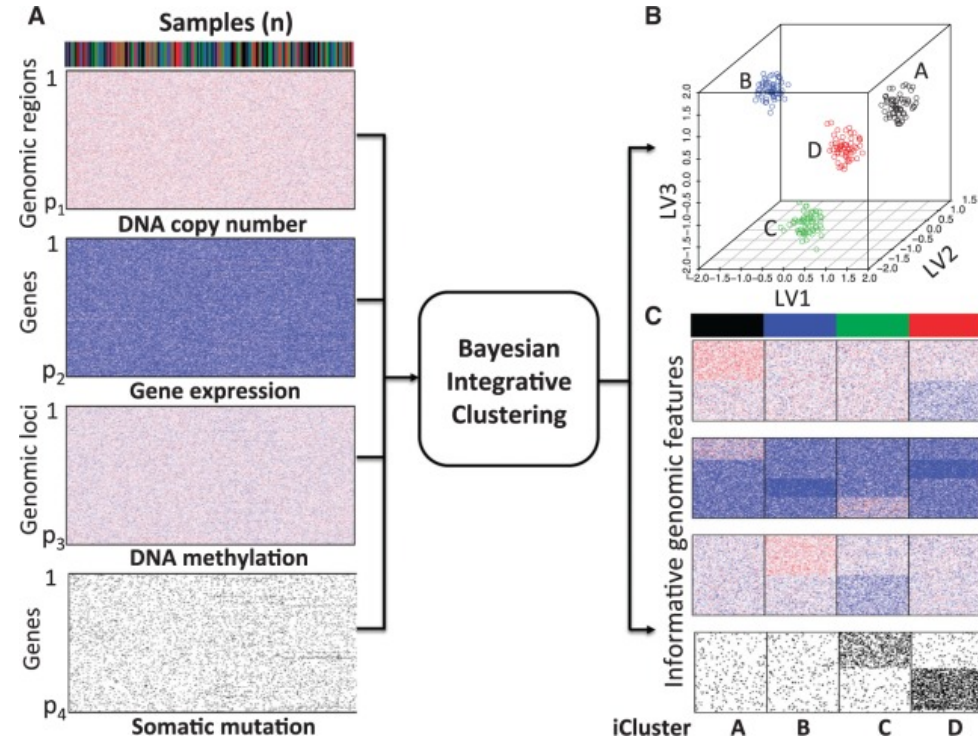
$\boldsymbol{\beta}_{jt} = (\beta_{0jt}, \beta_{1jt}, \dots, \beta_{kjt})^T$ *coefficients vector*

$\boldsymbol{\beta}_{jt} \sim \text{MVN}(\boldsymbol{\beta}_{0t}, \Sigma_{0t})$ *multivariate normal distribution*

$\sigma_{jt}^2 \sim \text{IG}(v_0/2, v_0 \sigma_0^2/2)$ *inverse-gamma distribution*

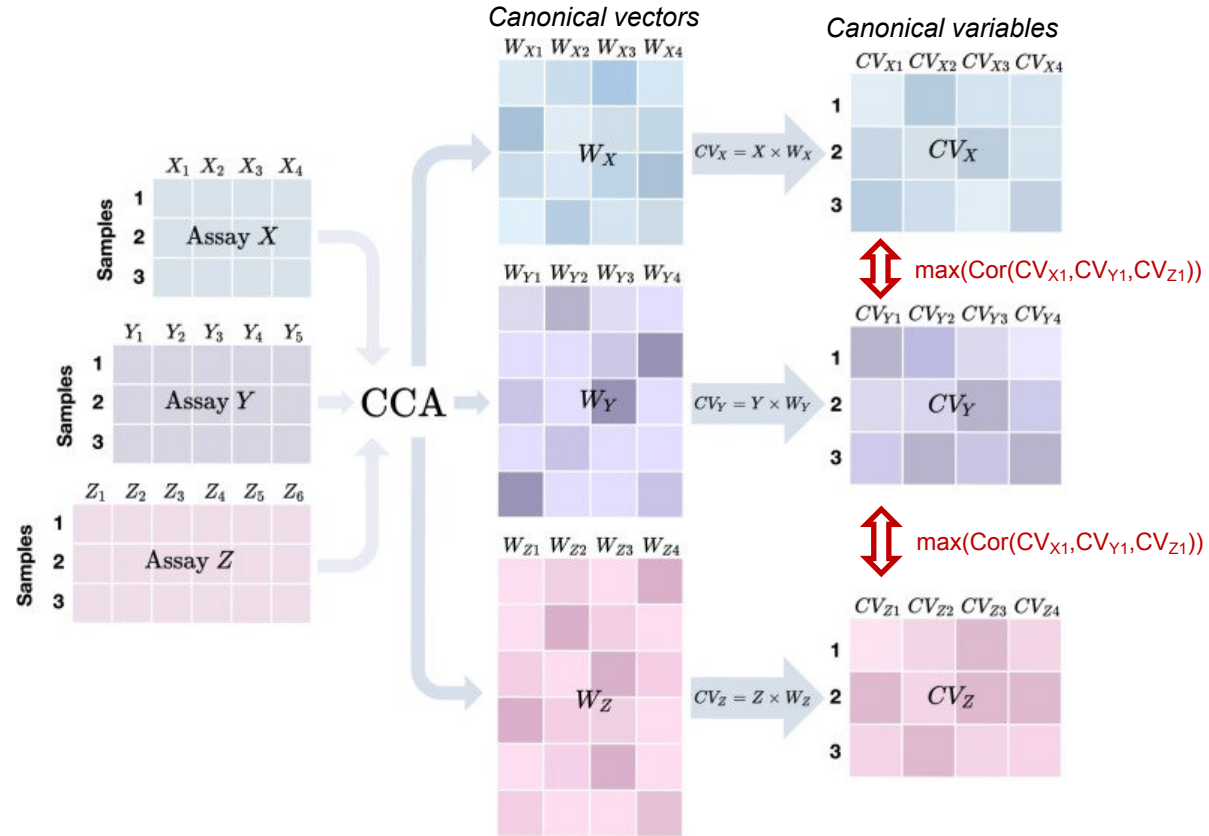
$\gamma_{jt} \sim \text{Bernoulli}(q_t)$. *Bernoulli distribution*

Use the Gibbs and Metropolis–Hasting sampling algorithms to sample for statistical inference.

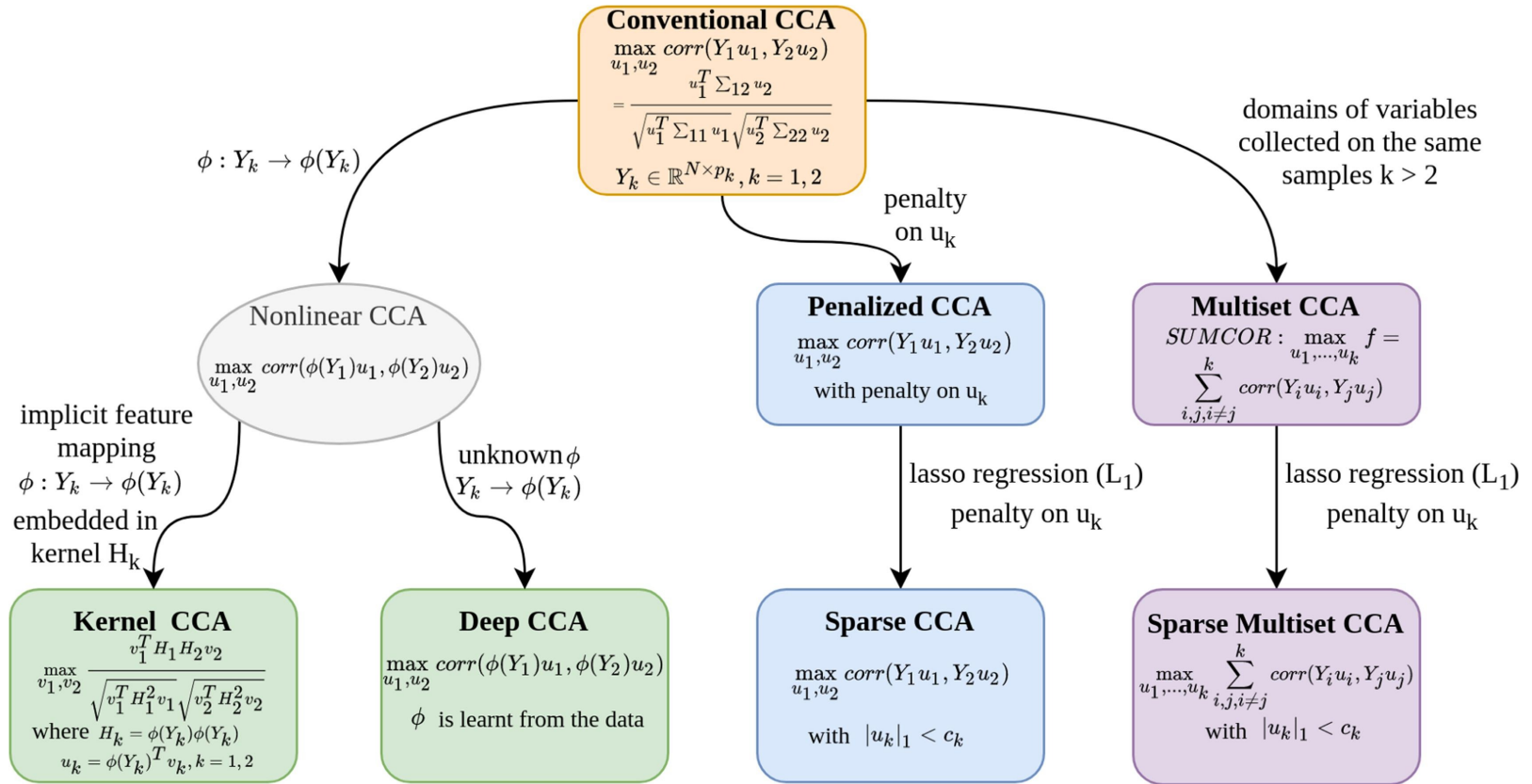




Canonical Correlation Analysis (CCA)

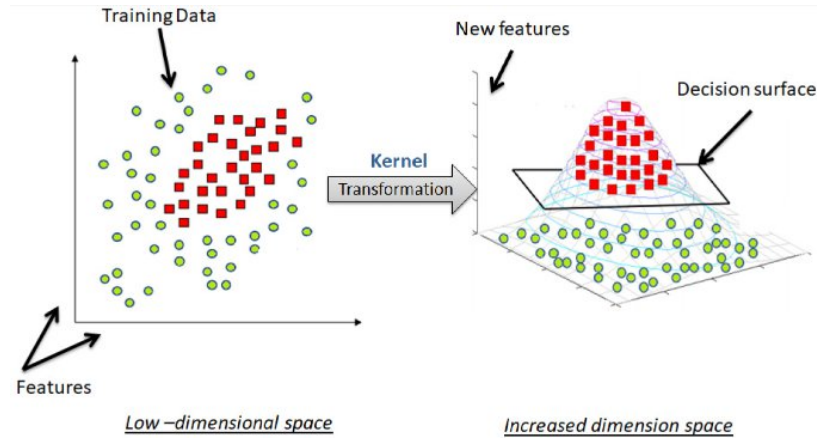


Jiang, MZ. et al. Canonical correlation analysis for multi-omics: Application to cross-cohort analysis. PLoS Genet. 2023



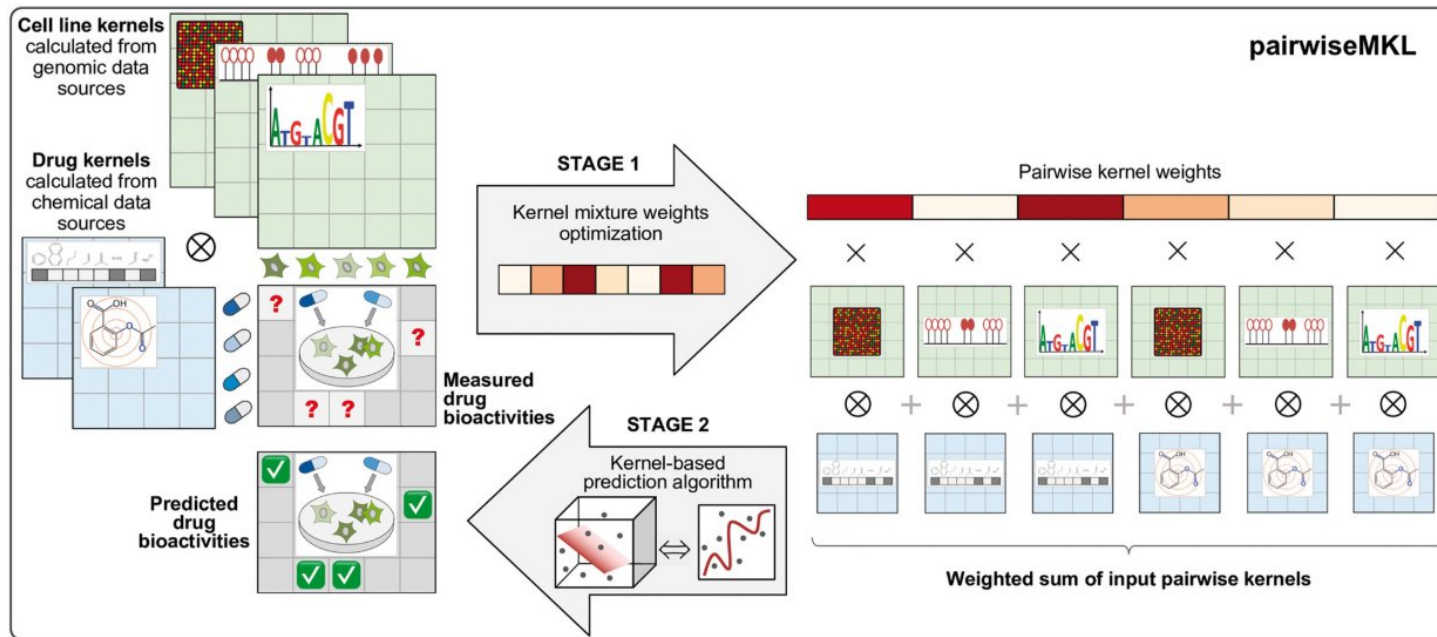


A kernel function maps input features onto a new space, which may be of higher dimension



Learn non-linear models from linear based algorithms using kernel functions $\forall(\mathbf{x}, \mathbf{x}') \in \mathcal{X}^2, K(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle_{\mathcal{X}}$

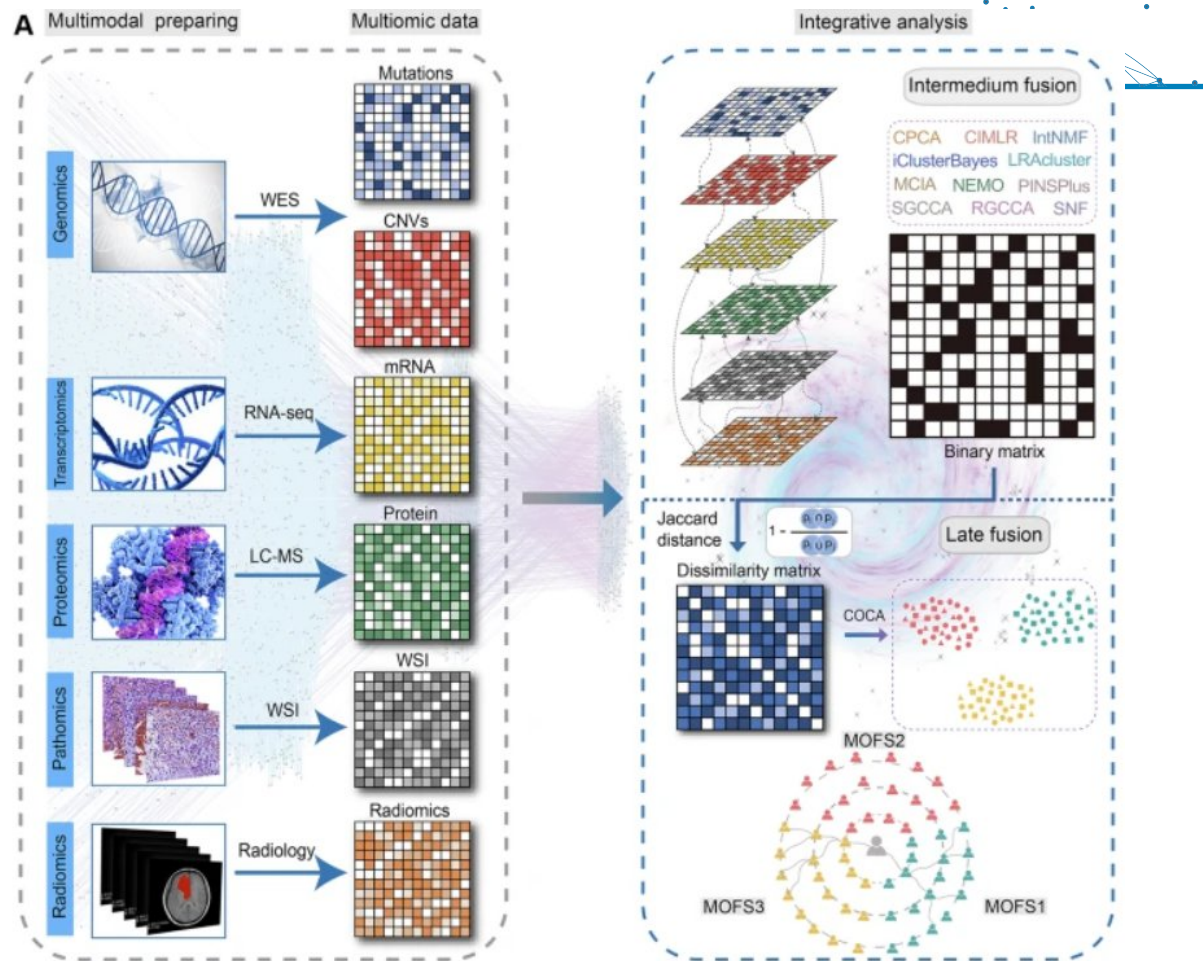
- Polynomial kernel : $\forall(\mathbf{x}, \mathbf{x}') \in \mathcal{X}^2, K(\mathbf{x}, \mathbf{x}') = (\langle \mathbf{x}, \mathbf{x}' \rangle_{\mathcal{X}})^{\gamma}$
- Gaussian kernel : $\forall(\mathbf{x}, \mathbf{x}') \in \mathcal{X}^2, K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2}\right)$



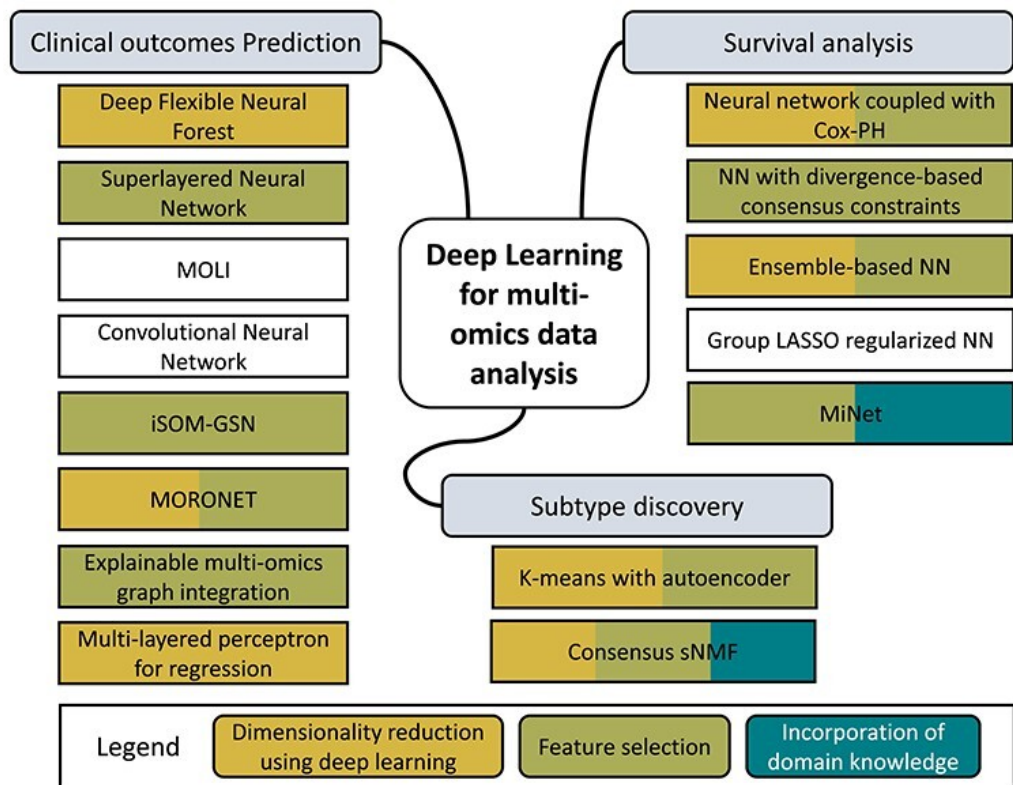
- 10 drugs kernels
- 12 cell-line kernels
- 312 proteins kernels
- 120 / 3120 pairwise kernels (Kronecker product)
- regularized model optimization to identify a subset of pairwise kernels weights

And combine them

- intermediate fusion of multimodal data through 11 algorithms
- each clustering results are converted into a binary matrix
- Jaccard index was calculated for the 11 matrices in order to assess the similarity between the samples
- late fusion of the results obtained from the 11 algorithms in order to derive the final clustering results using the COCA method (Clustering of Cluster Analysis) and arrive at consensus results



Liu, Z. et al. Multimodal fusion of radio-pathology and proteogenomics identify integrated glioma subtypes with prognostic and therapeutic opportunities. *Nat Commun* **16**, 3510 (2025).

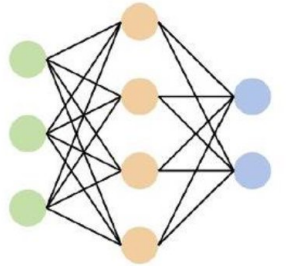


Kang, M., Ko, E., & Mersha, T. B. (2022). A roadmap for multi-omics data integration using deep learning. *Briefings in bioinformatics*, 23(1).



Non-generative methods

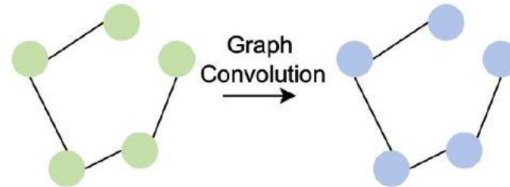
Feedforward Neural Network



- Inter-modality Interactions
- Biological Interpretability

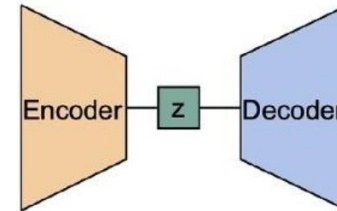
Input Layer Hidden Layer Output Layer

Graph Convolutional Neural Network



- Knowledge-guided Connectivity
- Data-Driven Connectivity

Autoencoder



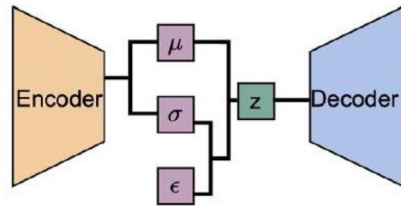
- Complementary Learning
- Consensus Learning
- Complementary and Consensus Learning
- Similarity Learning.

Ballard JL et al. Deep learning-based approaches for multi-omics data integration and analysis. *BioData Min.* 2024



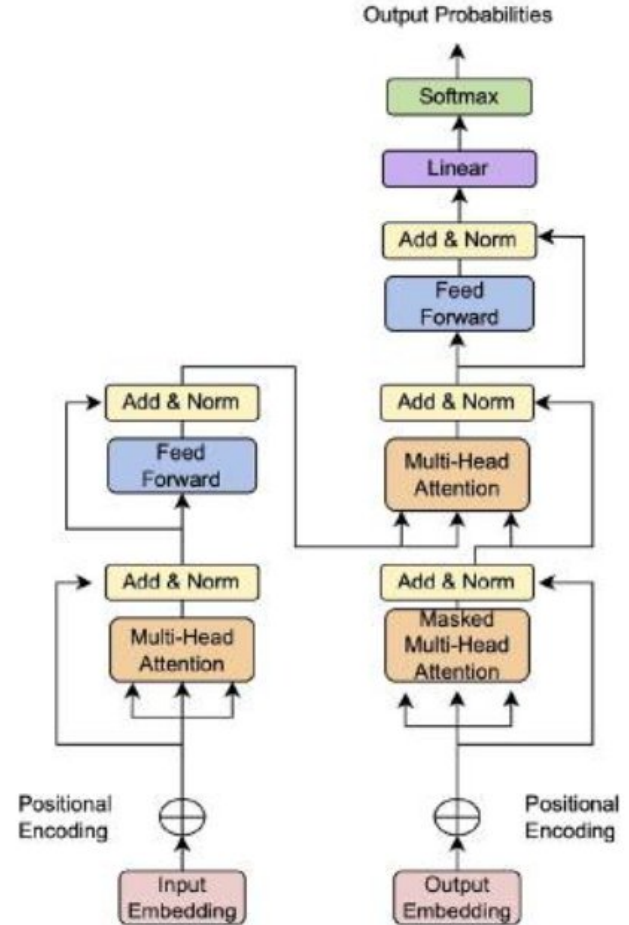
Generative methods

Variational Autoencoder

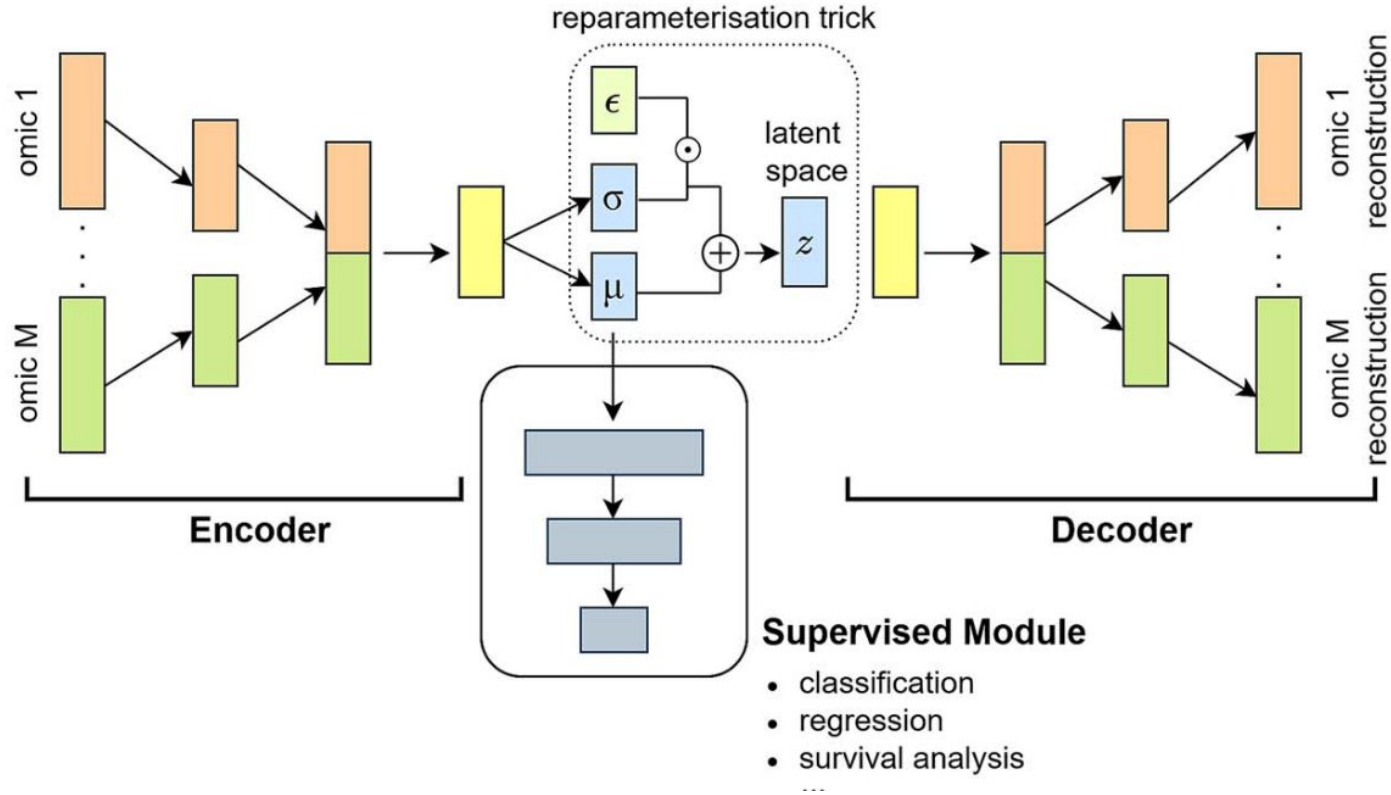


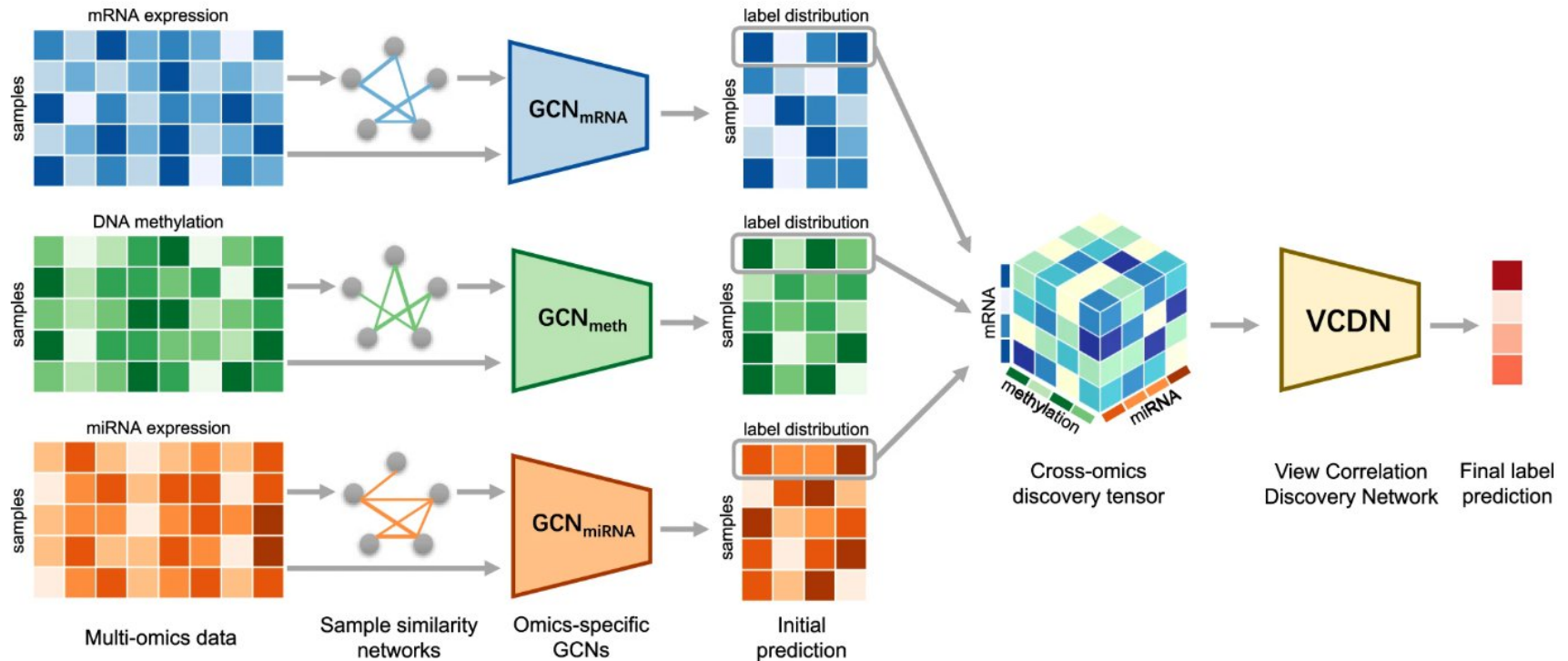
- Unsupervised learning for modeling inter-modality relationships
- Handling Unpaired Data
- Unsupervised learning for dimensionality reduction
- Supervised learning for generating task-relevant embeddings.

Generative Pretrained Transformer



Ballard JL et al. Deep learning-based approaches for multi-omics data integration and analysis. *BioData Min.* 2024

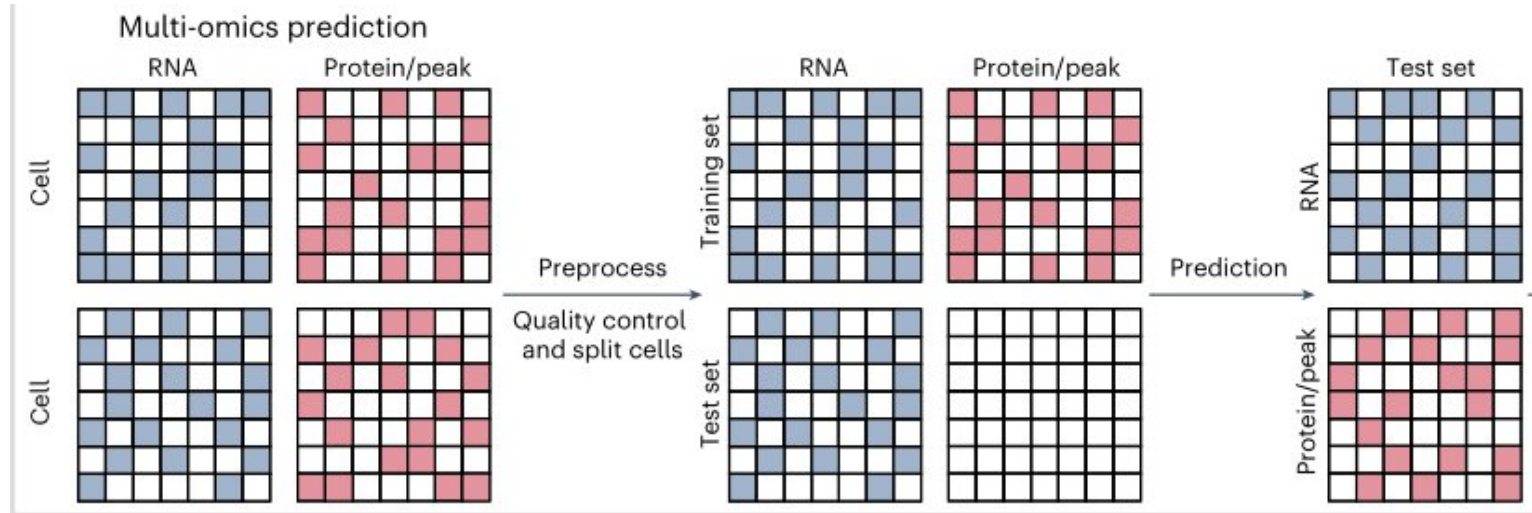




Wang, T. et al. MOGONET integrates multi-omics data using graph convolutional networks allowing patient classification and biomarker identification. Nat Commun 2021



- CITE-seq, REAP-seq ... : RNA expression and surface proteins abundance for a single cell;
- SHARE-seq, SNARE-seq ... : RNA expression and chromatin accessibility for a single cell;
- DOGMA-seq.... : RNA expression, chromatin accessibility and surface proteins abundance;
- ... (see Lim, J. et al. Advances in single-cell omics and multiomics for high-resolution molecular profiling. Exp Mol Med. 2024)



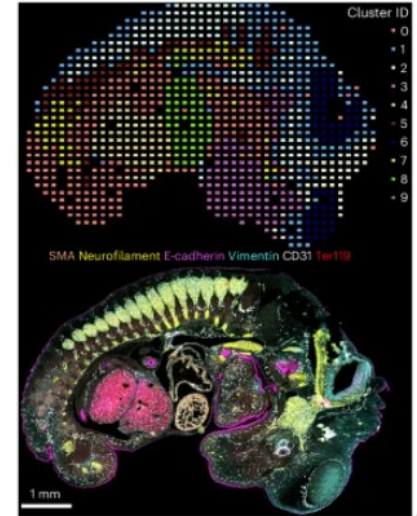
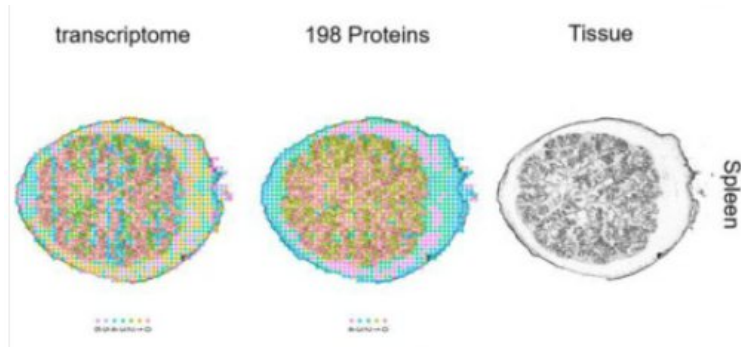


- From same single cell (matched)
 - Deep Learning: DCCA, DeepMAPS, scMVAE, totalVI
 - Probabilistic: MOFA+, BREM-SC, MIRA, MultiVelo
 - Graph / Metric-based: SCHEMA, Seurat v4, citeFUSE, FigR
- From different single cells (unmatched)
 - Deep Learning: GLUE, Cobolt
 - Probabilistic: MultiVI, LIGER
 - Graph: Spectrum, Seurat v5
 - Manifold Alignment: MMD-MA, UnionCom, Pamona
 - CCA: BindSC



- Spatial (only) transcriptomic: MERFISH, Stereo-seq
 - Prediction of non-spatial multi-omic data (e.g. proteins or open chromatin) by inferring the distribution of 'missing' omic data in space
 - Assume a common distribution between spatial and non-spatial omics data

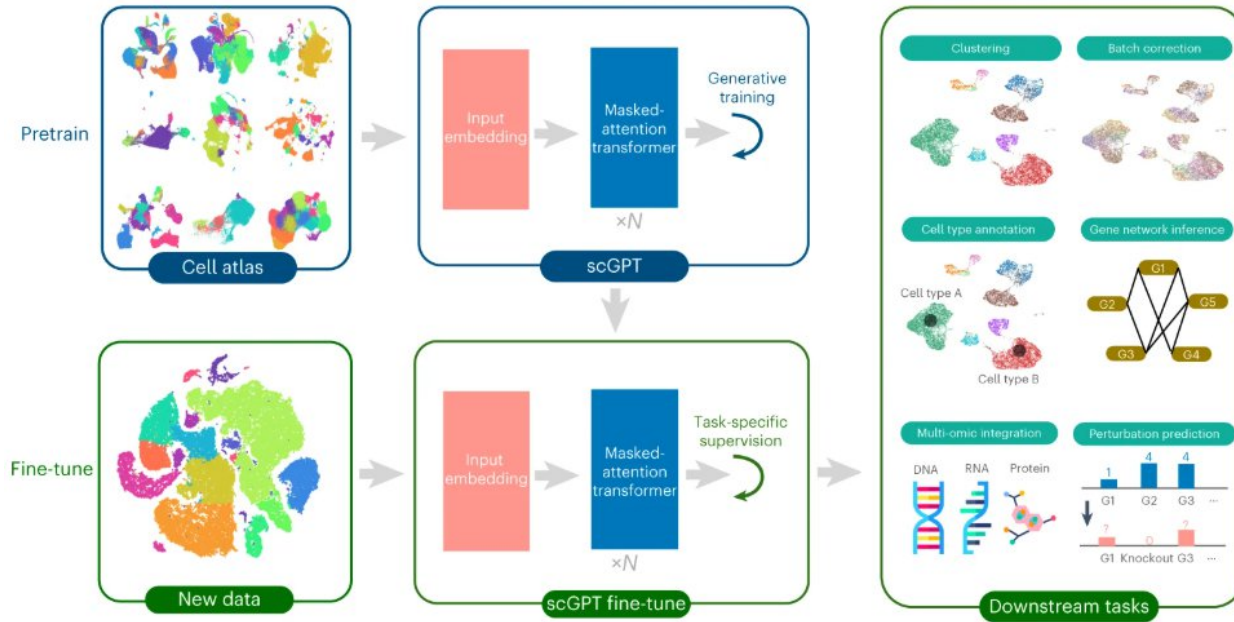
- Spatial multi-omics technology: spatial-CITE-seq and DBiT-seq



Liu, Y. et al. Spatial-CITE-seq: spatially resolved high-plex protein and whole transcriptome co-mapping. Res Sq 2022.
Enniful, A. et al. Integration of imaging-based and sequencing-based spatial omics mapping on the same tissue section via DBiTplus. Nat Methods 2026.



- Impute missing modalities or integrate multiple modalities
eg: *Geneformer*, *scGPT* and *scFoundation*



Two training steps:

- initial pretraining on large cell atlases (33 million non-disease human cells from 51 organs in the CELLxGENE collection)
- follow-up fine-tuning on smaller datasets for specific applications



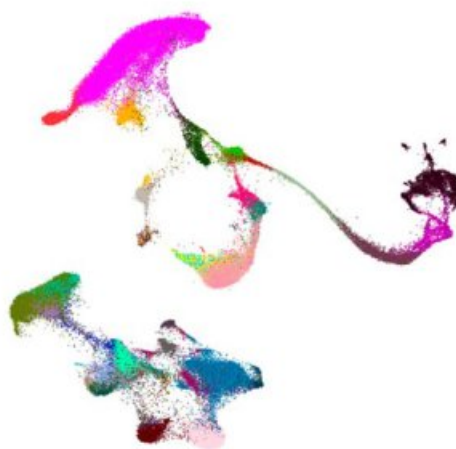
scGPT (fine-tuned)

Cell type, AvgBIO = 0.697



Seurat version 4

Cell type, AvgBIO = 0.600



- | | |
|---|--|
| ● B1 B IGKC ⁺ | ● CD8 ⁺ T naive |
| ● B1 B IGKC ⁻ | ● CD8 ⁺ T naive CD127 ⁺ CD26 ⁻ CD101 ⁻ |
| ● CD14 ⁺ mono | ● Erythroblast |
| ● CD16 ⁺ mono | ● G/M prog |
| ● CD4 ⁺ T CD314 ⁺ CD45RA ⁺ | ● HSC |
| ● CD4 ⁺ T activated | ● ILC |
| ● CD4 ⁺ T activated integrin β ₇ ⁺ | ● ILC1 |
| ● CD4 ⁺ T naive | ● Lymph prog |
| ● CD8 ⁺ T CD49f ⁺ | ● MAIT |
| ● CD8 ⁺ T CD57 ⁺ CD45RA ⁺ | ● MK/E prog |
| ● CD8 ⁺ T CD57 ⁺ CD45RO ⁺ | ● NK |
| ● CD8 ⁺ T CD69 ⁺ CD45RA ⁺ | ● NK CD158e1+ |
| ● CD8 ⁺ T CD69 ⁺ CD45RO ⁺ | ● Naive CD20 ⁺ B IGKC ⁺ |
| ● CD8 ⁺ T TIGIT ⁺ CD45RA ⁺ | ● Naive CD20 ⁺ B IGKC ⁻ |
| ● CD8 ⁺ T TIGIT ⁺ CD45RO ⁺ | ● Normoblast |

Fine-tuned scGPT model VS Seurat (v.4) on the CITE-seq BMDC dataset (paired RNA and protein data)

- 12 healthy human donors consisting of 12 batches (total of 90,261 cells)
- 13,953 genes and 134 surface proteins

Integration approaches : the good one ?

Integration approaches are not unic

- comparisons exist... for a given application)
- parametrization need expertise
- make your own comparisons
- keep an eye open

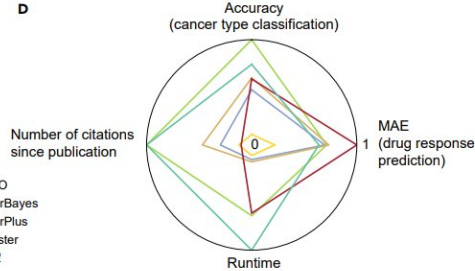
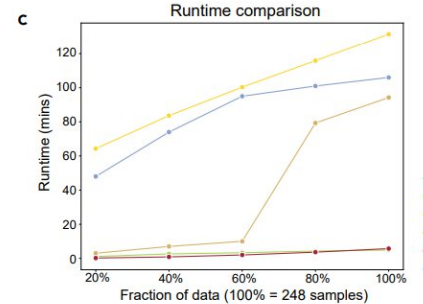
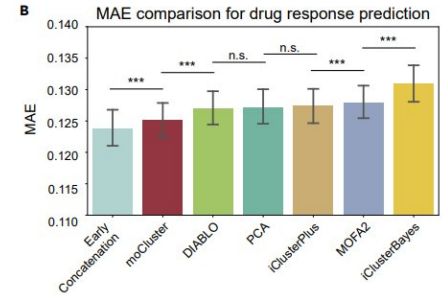
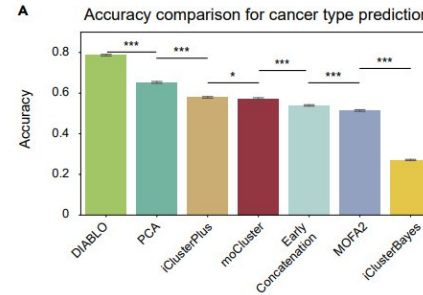
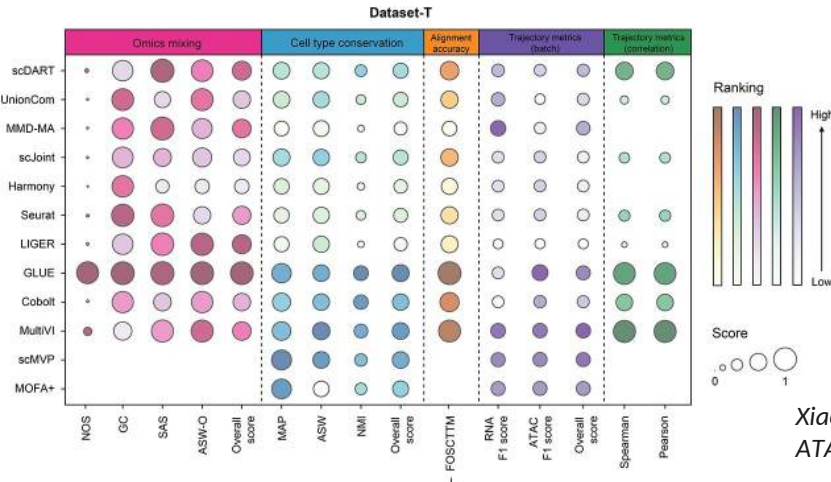


Figure 5. Benchmarking of machine learning-based integration tools using the CCLE multi-omics data

Cai, Z. et al. Machine learning for multi-omics data integration in cancer. *iScience*. 2022

Xiao, C. et al. Benchmarking multi-omics integration algorithms across single-cell RNA and ATAC data. *Brief Bioinform*. 2024



Integration approaches are not magic!

You will still need to:

- carefully check design and confounding factors
- perform specific data pre-processing for each omic (at least primary analysis)
- impute missing values* (different meaning → different strategy)
- choose your integration strategy based on your objective and your data (ex. matching between omics) → still no standard pipelines
- the best isn't necessarily the latest
- some omics generate more noise than answers



Table 1. Key portals for accessing publicly available multi-omics datasets

| Name | URL | Omic and other data types | Notes |
|--|---|---|---|
| TCGA (Campbell et al., 2020) | https://portal.gdc.cancer.gov/ | <ul style="list-style-type: none"> • Genomics • Epigenomics • Transcriptomics | <ul style="list-style-type: none"> • Tumor data • Large coverage of tumors |
| ICGC (Campbell et al., 2020) | https://dcc.icgc.org/ | <ul style="list-style-type: none"> • Genomics • Transcriptomics | <ul style="list-style-type: none"> • Tumor data • Powerful online analytics tools |
| CPTAC | https://cptac-data-portal.georgetown.edu/cptacPublic/ | <ul style="list-style-type: none"> • Proteomics | <ul style="list-style-type: none"> • Tumor data • The largest proteomic data portal |
| COSMIC Cell Lines (lorio et al., 2016) | https://cancer.sanger.ac.uk/cell_lines | <ul style="list-style-type: none"> • Genomics • Epigenomics • Transcriptomics • Drug response • CRISPR-Cas9 screen | <ul style="list-style-type: none"> • Cancer cell line data • Manually curated • Large coverage of cell lines |
| DepMap (Broad, 2020) | https://depmap.org/portal/ | <ul style="list-style-type: none"> • Genomics • Epigenomics • Transcriptomics • Proteomics • Drug response • CRISPR-Cas9 screen | <ul style="list-style-type: none"> • Cancer cell line data • Large coverage of omic types • Powerful online tools |
| COSMIC (Tate et al., 2019) | https://cancer.sanger.ac.uk/cosmic | <ul style="list-style-type: none"> • Genomics • Epigenomics • Transcriptomics | <ul style="list-style-type: none"> • Tumor data • Manually curated • Focus on genomics • Overlap with other portals |



Table 1. Key portals for accessing publicly available multi-omics datasets

| Name | URL |
|--|---|
| TCGA (Campbell et al., 2020) | https://portal.gdc.cancer.gov/ |
| ICGC (Campbell et al., 2020) | https://dcc.icgc.org/ |
| CPTAC | https://cptac-data-portal.georgetown.edu/cptacPublic/ |
| COSMIC Cell Lines (Iorio et al., 2016) | https://cancer.sanger.ac.uk/cell_lines |
| DepMap (Broad, 2020) | https://depmap.org/portal/ |
| COSMIC (Tate et al., 2019) | https://cancer.sanger.ac.uk/cosmic |

Resource | [Open access](#) | Published: 13 March 2025

Human BioMolecular Atlas Program (HuBMAP): 3D Human Reference Atlas construction and usage

[Katy Börner](#) ✉, [Philip D. Blood](#), [Jonathan C. Silverstein](#), [Matthew Ruffalo](#), [Rahul Satija](#), [Sarah A. Teichmann](#), [Gloria J. Pryhuber](#), [Ravi S. Misra](#), [Jeffrey M. Purkerson](#), [Jean Fan](#), [John W. Hickey](#), [Gesmira Molla](#), [Chuan Xu](#), [Yun Zhang](#), [Griffin M. Weber](#), [Yashvardhan Jain](#), [Danial Qaurooni](#), [Yongxin Kong](#), [HRA Team](#), [Andreas Bueckle](#) ✉ & [Bruce W. Herr II](#) ✉

Nature Methods **22**, 845–860 (2025) | [Cite this article](#)

- Transcriptomics
- Drug response
- CRISPR-Cas9 screen
- Genomics
- Epigenomics
- Transcriptomics
- Proteomics
- Drug response
- CRISPR-Cas9 screen
- Genomics
- Epigenomics
- Transcriptomics

JOURNAL ARTICLE

scMMO-atlas: a single cell multimodal omics atlas and portal for exploring fine cell heterogeneity and cell dynamics

[Wenwen Cheng](#), [Changhui Yin](#), [Shiya Yu](#), [Xi Chen](#), [Ni Hong](#), [Wenfei Jin](#) ✉

Nucleic Acids Research, Volume 53, Issue D1, 6 January 2025, Pages D1186–D1194, <https://doi.org/10.1093/nar/gkae821>

Published: 24 September 2024 **Article history** ▾

- Focus on genomics
- Overlap with other portals



Profiler Zirem, Y. et al. Profiler: an open web platform for multi-omics analysis, *Bioinformatics* 2026

iSODA Olivier-Jimenez, D. et al. iSODA: A Comprehensive Tool for Integrative Omics Data Analysis in Single- and Multi-Omics Experiments. *Anal Chem.* 2025

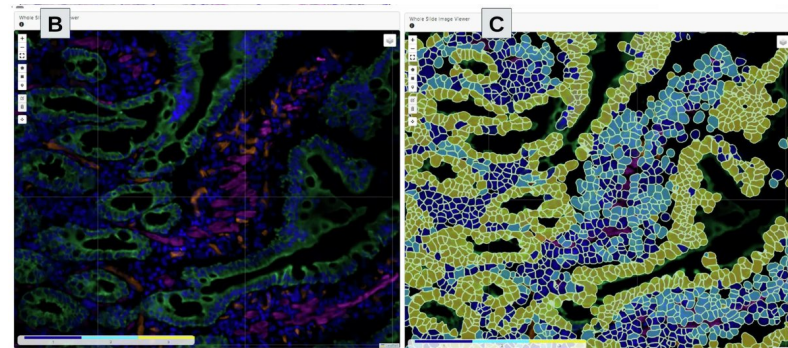
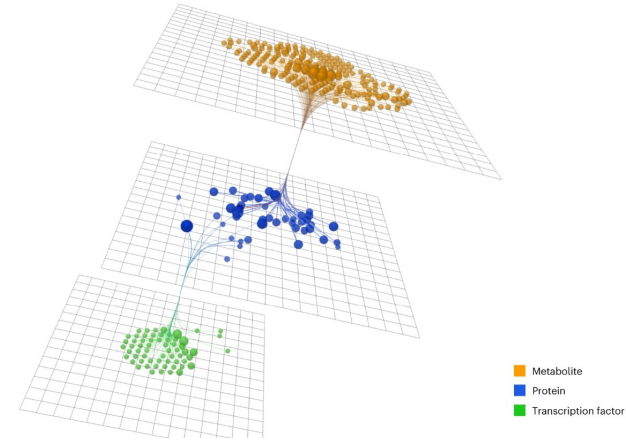
Analyst suite Ewald, J.D. et al. Web-based multi-omics integration using the Analyst software suite. *Nat Protoc.* 2024

FUSION Border, SP. et al. FUSION: a web-based application for in-depth exploration of multi-omics data with brightfield histology. *Nat Commun.* 2025

NOODAI Totu, T. et al. NOODAI: a webserver for network-oriented multi-omics data analysis and integration pipeline, *Bioinformatics* 2025

OmNI Potter, G. Et al. OmNI: a modular open-source framework for interactive multi-omics data integration and visualization. *NAR Genom Bioinform.* 2026

Galaxy-P project (Galaxy-P Project. galaxyp.org.)





Subramanian I, Verma S, Kumar S, Jere A, Anamika K. *Multi-omics Data Integration, Interpretation, and Its Application. Bioinform Biol Insights, 2020.*

Picard M. et al. *Integration strategies of multi-omics data for machine learning analysis. Comput Struct Biotechnol J. 2021.*

Baião AR. et al. *A technical review of multi-omics data integration methods: from classical statistical to deep generative approaches. Brief Bioinform. 2025*

Wróbel, S. et al. *Data integration through canonical correlation analysis and its application to OMICs research. J Biomed Inform. 2024*

Kang, M., Ko, E., & Mersha, T. B. *A roadmap for multi-omics data integration using deep learning. Briefings in bioinformatics, 2022.*

Xiao, C. et al. *Benchmarking multi-omics integration algorithms across single-cell RNA and ATAC data. Brief Bioinform. 2024*

