



SCIENCES ET
TECHNOLOGIES DE
L'INFORMATION ET DE
LA COMMUNICATION



Apprentissage profond pour l'analyse des organoïdes : modélisation par graphes des architectures cellulaires 3D

Alexandre Martin

Laboratoire d'Informatique, de Signaux et Systèmes de Sophia Antipolis (I3S)

UMR7271 UCA CNRS

Présentée en vue de l'obtention du grade de docteur en Informatique d'Université Côte d'Azur

Dirigée par : Xavier DESCOMBES, Directeur de recherche, I3S, INRIA, Université Côte

d'Azur

Soutenue le : 17 décembre 2025

Devant le jury, composé de :

Lionel FILLATRE, Professeur des universités, I3S, Université Côte d'Azur Alin ACHIM, Professeur, University of Bristol

Daniel RACOCEANU, Professeur, Sorbonne Université

Francesco PONZIO, Assistant Professeur, Politecnico di Torino Stéphan CLAVEL, Directeur de recherche,

IPMC, Université Côte d'Azur

APPRENTISSAGE PROFOND POUR L'ANALYSE DES ORGANOÏDES : MODÉLISATION PAR GRAPHES DES ARCHITECTURES CELLULAIRES 3D

Deep Learning for Organoid Analysis: Graph-Based Modeling of 3D Cellular Architectures

Alexandre MARTIN

 \bowtie

Jury:

Président du jury

Lionel FILLATRE, Professeur des universités, I3S, Université Côte d'Azur

Rapporteurs

Alin ACHIM, Professeur, University of Bristol Daniel RACOCEANU, Professeur, Sorbonne Université

Examinateurs

Francesco PONZIO, Assistant Professeur, Politecnico di Torino Stéphan CLAVEL, Directeur de recherche, IPMC, Université Côte d'Azur

Directeur de thèse

Xavier DESCOMBES, Directeur de recherche, I3S, INRIA, Université Côte d'Azur

Université Côte d'Azur	

Alexandre Martin Apprentissage profond pour l'analyse des organoïdes : modélisation par graphes des architectures cellulaires 3D xvi+56 p.

À toi lecteur <3 :*

Résumé

Les organoïdes, ces mini-organes cultivés *in vitro*, révolutionnent la recherche biomédicale en offrant des modèles tridimensionnels qui reproduisent la complexité des tissus humains. Cependant, leur analyse reste largement tributaire de méthodes manuelles, lentes et sujettes à des biais d'interprétation. Ces structures, composées de cellules organisées en réseaux d'interactions spatiales et fonctionnelles, nécessitent des outils capables de capturer non seulement leur morphologie, mais aussi les relations cellulaires qui déterminent leur mode de fonctionnement. C'est dans ce contexte que les réseaux de neurones sur graphes (Graph Neural Networks, GNN) émergent comme une solution particulièrement adaptée, permettant de modéliser les organoïdes non plus comme des images statiques, mais comme des systèmes relationnels où chaque cellule est un nœud connecté à ses voisines par des liens reflétant des interactions biologiques.

Cette thèse propose une approche innovante pour la modélisation et la classification automatisée des organoïdes à partir de graphes cellulaires, en exploitant pleinement le potentiel des GNN. Contrairement aux méthodes classiques basées sur des descriptions manuelles ou des réseaux de neurones convolutifs, qui analysent les images pixel par pixel, les GNN permettent d'intégrer des informations structurelles et contextuelles, en représentant chaque organoïde comme un réseau où les nœuds encodent des propriétés cellulaires (taille, forme, expression de marqueurs) et les arêtes capturent les relations spatiales. Cette représentation relationnelle ouvre la voie à une classification plus fine et plus interprétable, capable de distinguer des phénotypes subtils – comme des stades précoces de différenciation ou des altérations pathologiques – qui échappent aux approches traditionnelles.

Pour surmonter les défis posés par la rareté des données annotées et la variabilité intrinsèque des organoïdes, cette thèse développe une pipeline complète, depuis la construction de graphes cellulaires à partir d'images de microscopie jusqu'à l'apprentissage robuste de modèles de GNN. Une attention particulière est portée à la génération de données synthétiques via des modèles génératifs de graphes, afin d'enrichir les jeux d'entraînement et d'explorer des scénarios rares ou extrêmes. Enfin, des méthodes d'interprétation des prédictions sont mises en œuvre pour rendre les résultats exploitables par les biologistes, en identifiant les cellules et les interactions les plus déterminantes dans la classification.

Les applications de cette approche sont multiples : criblage à grande échelle de composés pharmaceutiques, diagnostic précoce de maladies à partir d'organoïdes dérivés de patients, ou encore optimisation des protocoles de culture pour standardiser la production d'organoïdes. À plus long terme, cette thèse jette les bases d'une analyse globale combinant imagerie, graphes cellulaires et données omiques, ouvrant la voie à une compréhension plus profonde des mécanismes biologiques sous-jacents et à des avancées en médecine personnalisée.

Abstract

Organoids—miniaturized, three-dimensional *in vitro* cultures that replicate the complexity of human tissues—are revolutionizing biomedical research. Yet their analysis remains heavily reliant on manual methods that are time-consuming, low-throughput, and prone to interpretative bias. These structures, composed of cells organized into spatial and functional interaction networks, demand analytical tools capable of capturing not only their morphology but also the cellular relationships that govern their behavior. In this context, Graph Neural Networks (GNNs) emerge as a particularly well-suited solution, enabling organoids to be modeled not as static images but as relational systems, where each cell is a node connected to its neighbors via edges representing biological interactions.

This thesis introduces an innovative framework for the automated modeling and classification of organoids using cellular graphs, fully leveraging the potential of GNNs. Unlike conventional approaches—based on manual descriptors or convolutional neural networks (CNNs), which analyze images pixel-by-pixel—GNNs integrate structural and contextual information by representing each organoid as a network. In this framework, nodes encode cellular properties (e.g., size, shape, marker expression) while edges capture spatial relationships. This relational representation enables finer and more interpretable classification, capable of distinguishing subtle phenotypes—such as early differentiation stages or pathological alterations—that elude traditional methods.

To address challenges posed by limited annotated data and the intrinsic variability of organoids, this work develops a comprehensive pipeline, from constructing cellular graphs from microscopy images to robust GNN training. Particular emphasis is placed on synthetic data generation via graph generative models to augment training sets and explore rare or extreme scenarios. Additionally, interpretability methods are implemented to make predictions actionable for biologists, highlighting the most decisive cells and interactions driving classification.

The applications of this approach are far-reaching: high-throughput drug screening, early disease diagnosis from patient-derived organoids, and optimization of culture protocols to standar-dize organoid production. In the long term, this thesis lays the groundwork for holistic multi-modal analysis—integrating imaging, cellular graphs, and omics data—to deepen our understanding of underlying biological mechanisms and advance precision medicine.

Remerciements

Merci!

Table des matières

1	Intr	oductio	n générale
	1.1	Contex	kte et motivation
		1.1.1	Organoïdes : révolution en biologie cellulaire et médecine régénérative .
		1.1.2	Applications thérapeutiques et criblage de médicaments
		1.1.3	Verrous scientifiques: quantification et standardisation
	1.2	Problé	matique scientifique
		1.2.1	Défis de l'analyse quantitative d'organoïdes 3D
		1.2.2	Limites des méthodes actuelles (manuelles, CNN 3D)
		1.2.3	Besoin d'approches structurelles adaptées
	1.3	Contri	butions de la thèse
		1.3.1	Pipeline automatisé de bout en bout pour organoïdes 3D
		1.3.2	Représentation par graphes géométriques et GNNs
		1.3.3	Génération de données synthétiques contrôlées
		1.3.4	Outils open-source pour la communauté
	1.4	Organi	isation du manuscrit
2	t		
	2.1	Organo	oïdes : biologie et applications
		2.1.1	Définitions et types d'organoïdes
		2.1.2	Mécanismes de formation et auto-organisation
		2.1.3	Applications : recherche fondamentale, drug screening, médecine person-
			nalisée
		2.1.4	Biomarqueurs et phénotypes d'intérêt
	2.2	Analys	se d'images biomédicales 3D
		2.2.1	Modalités d'imagerie
		2.2.2	Défis spécifiques : résolution, bruit, artefacts
		2.2.3	Contraintes computationnelles
	2.3	Métho	des d'analyse existantes
		2.3.1	Analyse manuelle : avantages et limites
		2.3.2	Segmentation cellulaire
		2.3.3	Approches par vision par ordinateur
		2.3.4	Méthodes basées graphes en histopathologie
	2.4	Positio	onnement de la thèse
		2.4.1	Lacunes identifiées dans la littérature
		2.4.2	Originalité de l'approche proposée
		2.43	Verrous scientifiques et techniques adressés

3	Fon	Condements théoriques					
	3.1	Théorie des graphes					
		3.1.1 Définitions formelles					
		3.1.2 Représentations matricielles					
		3.1.3 Métriques topologiques					
		3.1.4 Graphes géométriques vs abstraits					
	3.2	Graph Neural Networks: principes					
		3.2.1 Motivations et nécessité d'architectures spécialisées					
		3.2.2 Paradigme du message passing neural network (MPNN) 10					
		3.2.3 Couches de convolution sur graphes					
		3.2.4 Pooling et agrégation globale					
	3.3	Architectures GNN standards					
		3.3.1 Approches spectrales vs spatiales					
		3.3.2 Graph Convolutional Networks (GCN)					
		3.3.3 Graph Attention Networks (GAT)					
		3.3.4 GraphSAGE, GIN, et variantes modernes					
	3.4	GNNs géométriques : extensions E(3)-équivariantes					
		3.4.1 Symétries géométriques : translations, rotations, réflexions					
		3.4.2 Invariance vs équivariance : définitions et implications					
		3.4.3 Equivariant Graph Neural Networks (EGNN)					
		3.4.4 SchNet, DimeNet et architectures pour données 3D					
		3.4.5 Applications en chimie, biologie structurale, physique					
	3.5	Expressivité et limitations théoriques					
		3.5.1 Weisfeiler-Lehman test et pouvoir de discrimination					
		3.5.2 Over-smoothing et profondeur des réseaux					
		3.5.3 Over-squashing et propagation d'information					
	3.6	Statistiques spatiales pour données ponctuelles					
		3.6.1 Processus ponctuels : définitions et propriétés					
		3.6.2 Fonctions de Ripley (K, F, G) sur surfaces courbes					
		3.6.3 Tests d'agrégation et de régularité spatiale					
		3.6.4 Applications en biologie cellulaire					
4	Mét	hodologie et pipeline de traitement 15					
•	4.1	Architecture générale du pipeline					
	7.1	4.1.1 Vue d'ensemble : de l'image brute au graphe					
		4.1.2 Choix de conception et compromis					
		4.1.3 Justifications scientifiques et techniques					
	4.2	Acquisition et prétraitement des images					
	7.2	4.2.1 Protocoles expérimentaux et marqueurs fluorescents					
		4.2.2 Normalisation d'intensité et débruitage					
		4.2.3 Correction d'artefacts d'acquisition					
		4.2.4 Gestion de la variabilité inter-expérimentale					
	4.3	Segmentation cellulaire automatisée					
	+.3	4.3.1 Revue des méthodes					
		4.3.2 Comparaison quantitative					
		4.3.3 Choix et paramétrage de l'outil optimal					
		T.S.S CHOIX OF PARAMORAGE WE I DUM OPHIMAL					

		4.3.4 Validation qualitative avec experts biologistes
	4.4	Extraction et séparation des organoïdes
		4.4.1 Conversion en nuages de points 3D
		4.4.2 Clustering spatial
		4.4.3 Filtrage des artefacts et débris cellulaires
		4.4.4 Identification et isolation des organoïdes individuels
	4.5	Construction de graphes géométriques
		4.5.1 Définition des nœuds : propriétés cellulaires
		4.5.2 Stratégies de connectivité
		4.5.3 Normalisation et standardisation des features
		4.5.4 Analyse de sensibilité aux hyperparamètres
	4.6	Génération de données synthétiques
		4.6.1 Motivation : rareté et coût des annotations
		4.6.2 Processus ponctuels sur la sphère
		4.6.3 Diagrammes de Voronoï 3D pour géométrie réaliste
		4.6.4 Simulation de phénotypes contrôlés
		4.6.5 Validation statistique
		4.6.6 Augmentation par perturbations contrôlées
	4.7	Architectures GNN implémentées
		4.7.1 Choix d'architectures adaptées aux graphes géométriques
		4.7.2 Baseline : GCN, GAT pour ablation
		4.7.3 EGNN pour équivariance E(3)
		4.7.4 Adaptations et modifications proposées
	4.8	Entraînement et optimisation
		4.8.1 Fonctions de perte pour classification
		4.8.2 Stratégies de régularisation
		4.8.3 Hyperparamètres et recherche automatique
		4.8.4 Validation croisée adaptée aux petits datasets
5	Exp	érimentations et résultats 21
	5.1	Protocole expérimental
		5.1.1 Datasets: description et caractéristiques
		5.1.2 Métriques d'évaluation
		5.1.3 Conditions expérimentales et reproductibilité
	5.2	Validation des données synthétiques
		5.2.1 Analyse des statistiques spatiales
		5.2.2 Réalisme géométrique et morphologique
		5.2.3 Diversité des phénotypes générés
		5.2.4 Pertinence pour le pré-entraînement
	5.3	Résultats sur données synthétiques
		5.3.1 Capacité de discrimination entre processus ponctuels
		5.3.2 Comparaison des architectures GNN
		5.3.3 Études d'ablation
		5.3.4 Analyse de sensibilité aux hyperparamètres
	5.4	Résultats sur données réelles
		5.4.1 Classification de phénotypes biologiques

		5.4.2	Comparaison avec méthodes de référence
		5.4.3	Performances en fonction de la taille du dataset
		5.4.4	Généralisation inter-expérimentale
	5.5	Appro	che hybride : pré-entraînement + fine-tuning
		5.5.1	Pré-entraînement sur données synthétiques
		5.5.2	Fine-tuning sur données réelles
		5.5.3	Gains de performances et data efficiency
		5.5.4	Analyse des représentations apprises
	5.6	Interp	étabilité et analyse biologique
		5.6.1	Attention maps et cellules importantes
		5.6.2	Identification de patterns spatiaux discriminants
		5.6.3	Corrélation avec biomarqueurs connus
		5.6.4	Validation par experts biologistes
	5.7	Discus	ssion
		5.7.1	Forces de l'approche proposée
		5.7.2	Limitations et cas d'échec
		5.7.3	Comparaison critique avec l'état de l'art
		5.7.4	Compromis précision-interprétabilité-efficacité
6	Con	clusion	et perspectives 27
	6.1	Synthe	ese des contributions
		6.1.1	Récapitulatif des verrous levés
		6.1.2	Avancées méthodologiques
		6.1.3	Résultats expérimentaux majeurs
		6.1.4	Apports pour la communauté scientifique
	6.2	Limita	tions et défis
		6.2.1	Généralisabilité à d'autres types d'organoïdes
		6.2.2	Scalabilité aux très grands volumes
		6.2.3	Robustesse aux variations d'acquisition
		6.2.4	Nécessité d'annotations expertes
	6.3	Perspe	ctives à court terme
		6.3.1	Extension à d'autres phénotypes et pathologies
		6.3.2	Intégration de données multi-modales
		6.3.3	Amélioration de l'interprétabilité
		6.3.4	Validation clinique et transfert technologique
	6.4	Perspe	ctives à long terme
		6.4.1	Analyse spatio-temporelle: suivi longitudinal
		6.4.2	Modèles génératifs de graphes d'organoïdes
		6.4.3	Prédiction de réponse thérapeutique
		6.4.4	Vers une analyse holistique multi-échelles
	6.5		t scientifique et sociétal
		6.5.1	Accélération de la recherche sur organoïdes
		6.5.2	Applications en médecine personnalisée
		6.5.3	Réduction des coûts et du temps d'analyse
		654	Contribution aux alternatives à l'expérimentation animale 30

TA	BLE	DES MATIÈRES	xv
No	otation	ns	33
Ré	féren	nces	35
Lis	ste de	es figures	37
Lis	ste de	es définitions	39
Lis	ste de	es exemples	41
		Annexes	
A	Fond	damentaux du Deep Learning	47
	A. 1	Histoire et évolutions majeures	47
	A.2	Architectures classiques	47
		A.2.1 Perceptrons multicouches (MLP)	47
		A.2.2 Réseaux de neurones convolutifs (CNN)	47
		A.2.3 Réseaux de neurones récurrents (RNN, LSTM, GRU)	47
		A.2.4 Transformers	47
	A.3	Techniques d'optimisation	47
		A.3.1 Algorithmes d'optimisation	47
	A 1	A.3.2 Bonnes pratiques	48
	A.4	Régularisation	48 48
		A.4.1 Dropout	48
		A.4.3 Early stopping	48
	A.5	Bonnes pratiques d'entraînement	48
В	Com	apléments sur les graphes et GNNs	49
	B.1	Types de graphes exotiques	49
		B.1.1 Hypergraphes	49
		B.1.2 Graphes dynamiques	49
		B.1.3 Graphes hétérogènes	49
	B.2	Geometric Deep Learning : formalisme unifié	49
	B.3	Topological Data Learning	49
		B.3.1 Persistent homology	49
		B.3.2 Applications aux graphes	49
	B.4	Expressivité théorique : au-delà du WL-test	50
		B.4.1 Hiérarchie WL	50
	D 7	B.4.2 Higher-order GNNs	50
	B.5	Message passing généralisé	50

C	Thé	forie des processus ponctuels	51
	C .1	Processus de Poisson : propriétés et simulation	51
		C.1.1 Définition	51
		C.1.2 Propriétés	51
		C.1.3 Simulation	51
	C.2	Processus de Cox et processus log-gaussiens	51
	C.3	Processus de Gibbs et modèles énergétiques	51
	C.4	Estimation statistique et inférence	52
		C.4.1 Maximum de vraisemblance	52
		C.4.2 Méthodes basées simulation	52
	C.5	Processus ponctuels sur variétés	52
		C.5.1 Extension aux sphères	52
		C.5.2 Processus sur surfaces courbes	52
D	Déta	ails d'implémentation	53
	D.1	Technologies et bibliothèques	53
		D.1.1 Environnement logiciel	53
		D.1.2 Outils de segmentation	53
		D.1.3 Visualisation et analyse	53
	D.2	Architecture logicielle du pipeline	53
	D.3	Gestion des ressources computationnelles	54
		D.3.1 Hardware utilisé	54
		D.3.2 Optimisations	54
	D.4	Documentation des hyperparamètres	54
	D.5	Reproductibilité	54
E	Don	nnées et benchmarks	55
	E.1	Description détaillée des datasets	55
		E.1.1 Dataset synthétique	55
		E.1.2 Dataset réel	
	E.2	Protocoles d'annotation	55
	E.3	Statistiques descriptives complètes	55
	E.4	Accès aux données et code	
		E.4.1 Dépôt de code	
		E.4.2 Données	56
		E.4.3 Environnement reproductible	56

CHAPITRE 1

Introduction générale

1.1 Contexte et motivation

1.1.1 Organoïdes : révolution en biologie cellulaire et médecine régénérative

Les organoïdes, ces structures tridimensionnelles cultivées *in vitro* qui miment la complexité architecturale et fonctionnelle des organes humains, représentent une avancée majeure en biologie cellulaire et en médecine régénérative. Contrairement aux cultures cellulaires bidimensionnelles traditionnelles, les organoïdes reproduisent l'organisation spatiale, les interactions cellulaires et les gradients biochimiques caractéristiques des tissus *in vivo*.

1.1.2 Applications thérapeutiques et criblage de médicaments

Les applications des organoïdes s'étendent de la recherche fondamentale au développement thérapeutique. En oncologie, les organoïdes dérivés de patients permettent de tester *ex vivo* l'efficacité de traitements personnalisés. Dans le domaine des maladies génétiques, ils offrent des modèles cellulaires portant les mutations d'intérêt. Le criblage à haut débit de composés pharmaceutiques sur organoïdes promet d'accélérer la découverte de nouveaux médicaments tout en réduisant le recours à l'expérimentation animale.

1.1.3 Verrous scientifiques : quantification et standardisation

Malgré leur potentiel, l'exploitation optimale des organoïdes se heurte à des défis majeurs de quantification et de standardisation. L'analyse de leur morphologie tridimensionnelle complexe, de leur hétérogénéité intrinsèque, et l'identification de biomarqueurs phénotypiques requièrent des outils d'analyse avancés, actuellement insuffisants.

1.2 Problématique scientifique

1.2.1 Défis de l'analyse quantitative d'organoïdes 3D

L'analyse quantitative des organoïdes 3D présente plusieurs défis majeurs :

- Complexité morphologique : Les organoïdes présentent des architectures tridimensionnelles avec des arrangements cellulaires complexes difficiles à caractériser avec les méthodes traditionnelles.
- Hétérogénéité: Une variabilité importante existe tant au sein d'un même organoïde (hétérogénéité intra-organoïde) qu'entre différents organoïdes (hétérogénéité inter-organoïdes).

- **Contraintes computationnelles**: Les images 3D d'organoïdes peuvent atteindre plusieurs gigaoctets, posant des défis de stockage, de traitement et d'analyse.
- Absence de vérité terrain : Le manque de datasets annotés publics et la difficulté d'obtenir des annotations expertes fiables limitent le développement de méthodes d'apprentissage automatique.

1.2.2 Limites des méthodes actuelles (manuelles, CNN 3D)

Les approches actuelles d'analyse d'organoïdes présentent des limitations importantes :

- Analyse manuelle : Chronophage, sujette à la subjectivité et à la variabilité interobservateur, elle ne permet pas le passage à l'échelle nécessaire pour les études à haut débit.
- Réseaux de neurones convolutifs 3D : Bien que performants pour l'analyse d'images, les CNN 3D sont gourmands en mémoire et ne capturent pas efficacement les relations structurelles et topologiques entre cellules.
- Descripteurs manuels : Les approches basées sur des descripteurs de texture handcrafted manquent de généralité et nécessitent une expertise domaine importante.

1.2.3 Besoin d'approches structurelles adaptées

Face à ces limitations, il apparaît nécessaire de développer des approches qui exploitent explicitement la structure relationnelle des organoïdes. Les cellules au sein d'un organoïde forment un réseau d'interactions spatiales et fonctionnelles qu'il convient de modéliser directement, plutôt que de traiter l'organoïde comme une simple image tridimensionnelle.

1.3 Contributions de la thèse

1.3.1 Pipeline automatisé de bout en bout pour organoïdes 3D

Cette thèse propose un pipeline complet et automatisé pour l'analyse d'organoïdes 3D, depuis l'acquisition d'images de microscopie jusqu'à la classification de phénotypes. Ce pipeline intègre des étapes de segmentation cellulaire, d'extraction de features géométriques, de construction de graphes et de classification par Graph Neural Networks.

1.3.2 Représentation par graphes géométriques et GNNs

Une contribution majeure de ce travail réside dans la représentation des organoïdes sous forme de graphes géométriques, où chaque cellule constitue un nœud enrichi de propriétés morphologiques et d'intensités de marqueurs, et où les arêtes encodent les relations de voisinage spatial. Cette représentation permet l'application de Graph Neural Networks (GNNs) équivariants, capables de capturer les patterns structurels discriminants tout en respectant les symétries géométriques naturelles.

1.3.3 Génération de données synthétiques contrôlées

Pour pallier le manque de données annotées, nous proposons une approche de génération de données synthétiques basée sur la théorie des processus ponctuels. En simulant différents proces-

sus spatiaux (Poisson, Matérn, Strauss) sur des géométries sphériques et en construisant des diagrammes de Voronoï 3D, nous générons des organoïdes synthétiques aux propriétés statistiques contrôlées, permettant un pré-entraînement efficace des modèles.

1.3.4 Outils open-source pour la communauté

L'ensemble des outils développés, du pipeline de traitement aux architectures GNN, sera mis à disposition de la communauté scientifique sous licence open-source, accompagné de documentation et de tutoriels, afin de faciliter la reproductibilité et l'adoption par d'autres équipes de recherche.

1.4 Organisation du manuscrit

Ce manuscrit est organisé en six chapitres principaux complétés par cinq annexes techniques.

Le **Chapitre 2** présente un état de l'art exhaustif couvrant la biologie des organoïdes, les techniques d'imagerie 3D, les méthodes existantes d'analyse d'images biomédicales, et positionne nos contributions par rapport à la littérature.

Le **Chapitre 3** établit les fondements théoriques nécessaires à la compréhension de notre approche : théorie des graphes, Graph Neural Networks standard et géométriques, et statistiques spatiales pour processus ponctuels.

Le **Chapitre 4** décrit en détail la méthodologie proposée, depuis l'acquisition et le prétraitement des images jusqu'à la construction de graphes et l'architecture des GNNs, en passant par la génération de données synthétiques.

Le **Chapitre 5** présente les résultats expérimentaux obtenus sur données synthétiques et réelles, incluant des études d'ablation, des comparaisons avec l'état de l'art, et une analyse approfondie des performances et de l'interprétabilité.

Le **Chapitre 6** conclut en synthétisant les contributions, en discutant les limitations, et en proposant des perspectives à court et long terme pour ce domaine de recherche.

Les **annexes** fournissent des compléments théoriques sur le deep learning (Annexe A), les graphes et GNNs (Annexe B), la théorie des processus ponctuels (Annexe C), les détails d'implémentation (Annexe D), et la documentation des données (Annexe E).

Chapitre 2

État de l'art

2.1 Organoïdes: biologie et applications

2.1.1 Définitions et types d'organoïdes

Les organoïdes sont des structures tridimensionnelles auto-organisées cultivées *in vitro* à partir de cellules souches ou de tissus primaires. Différents types d'organoïdes ont été développés, reproduisant la structure et la fonction d'organes variés : organoïdes intestinaux, cérébraux, rénaux, hépatiques, pulmonaires, pancréatiques, etc. Chaque type présente des caractéristiques morphologiques et cellulaires spécifiques.

2.1.2 Mécanismes de formation et auto-organisation

La formation d'organoïdes repose sur les capacités d'auto-organisation cellulaire. En présence d'une matrice extracellulaire appropriée et de facteurs de croissance spécifiques, les cellules se différencient, prolifèrent et s'organisent spontanément selon des principes morphogénétiques similaires à ceux du développement embryonnaire.

2.1.3 Applications : recherche fondamentale, drug screening, médecine personnalisée

Les organoïdes trouvent des applications multiples en recherche biomédicale :

- Modélisation du développement : Étude des processus morphogénétiques et de différenciation
- Modélisation de maladies : Reproduction in vitro de pathologies génétiques ou infectieuses
- Criblage pharmacologique : Test à haut débit de l'efficacité et de la toxicité de composés
- Médecine personnalisée : Prédiction de réponse thérapeutique à partir d'organoïdes dérivés de patients
- Médecine régénérative : Source potentielle de tissus pour transplantation

2.1.4 Biomarqueurs et phénotypes d'intérêt

L'identification de phénotypes dans les organoïdes repose sur divers biomarqueurs : marqueurs de prolifération (Ki67), de mort cellulaire (caspase-3), de différenciation (marqueurs lignage-spécifiques), et de fonctionnalité. La caractérisation quantitative de ces phénotypes est cruciale pour exploiter pleinement le potentiel des organoïdes.

2.2 Analyse d'images biomédicales 3D

2.2.1 Modalités d'imagerie

Plusieurs modalités d'imagerie permettent d'observer les organoïdes en trois dimensions :

- **Microscopie confocale**: Haute résolution, mais acquisition lente et photoblanchiment
- Light-sheet microscopy : Imagerie rapide de grands volumes avec faible phototoxicité
- Microscopie multiphoton : Imagerie en profondeur avec réduction des dommages cellulaires
- Microscopie à expansion : Augmentation physique de la taille des échantillons pour améliorer la résolution

2.2.2 Défis spécifiques : résolution, bruit, artefacts

L'imagerie 3D d'organoïdes présente des défis techniques importants : résolution spatiale limitée en profondeur, diffusion de la lumière, hétérogénéité d'illumination, bruit de fond, aberrations optiques. Ces artefacts compliquent l'analyse automatisée et nécessitent des stratégies de prétraitement sophistiquées.

2.2.3 Contraintes computationnelles

Un organoïde imagé en haute résolution peut générer des volumes de données dépassant 2 Go par échantillon. Le traitement de telles images nécessite des ressources computationnelles importantes et des stratégies d'optimisation mémoire, limitant l'application de certaines approches de deep learning.

2.3 Méthodes d'analyse existantes

2.3.1 Analyse manuelle : avantages et limites

L'analyse manuelle par des experts reste la référence (*gold standard*) pour l'évaluation d'organoïdes. Elle permet une interprétation riche en contexte biologique mais souffre de limitations majeures : temps d'analyse élevé (plusieurs minutes par organoïde), subjectivité, variabilité interet intra-observateur, impossibilité de passage à l'échelle pour des études à haut débit.

2.3.2 Segmentation cellulaire

La segmentation automatique des cellules constitue une étape critique du pipeline d'analyse. Plusieurs approches ont été proposées :

- Méthodes géométriques : Détection de formes ellipsoïdales, limitée aux morphologies simples
- Watershed: Segmentation par bassins versants, sensible au sur-segmentation
- **Stardist**: Détection par prédiction de distances radiales, efficace pour noyaux convexes
- Cellpose : Architecture de deep learning avec champs de gradients, état de l'art actuel

2.3.3 Approches par vision par ordinateur

Les approches classiques de vision par ordinateur appliquées aux organoïdes incluent :

- Descripteurs de texture : Matrices de co-occurrence (Haralick), Local Binary Patterns, filtres de Gabor. Ces descripteurs manuellement conçus manquent de flexibilité et nécessitent une expertise domaine.
- CNN 2D : Analyse slice-by-slice, perte d'information 3D et cohérence spatiale
- CNN 3D: Limitations mémoire importantes, nécessitent downsampling massif, sensibles aux variations d'acquisition (luminosité, contraste)

2.3.4 Méthodes basées graphes en histopathologie

Quelques travaux ont exploré l'utilisation de graphes pour l'analyse de tissus en histopathologie 2D, représentant les cellules comme des nœuds. Ces approches ont montré des résultats prometteurs pour la classification de cancers, mais leur extension aux organoïdes 3D n'a pas été explorée.

2.4 Positionnement de la thèse

2.4.1 Lacunes identifiées dans la littérature

À notre connaissance, aucune méthode automatisée n'exploite pleinement la structure relationnelle 3D des organoïdes via des Graph Neural Networks géométriques. Les approches existantes traitent les organoïdes comme des images, sans modéliser explicitement le réseau d'interactions cellulaires qui gouverne leur comportement.

2.4.2 Originalité de l'approche proposée

Notre approche se distingue par :

- La représentation des organoïdes comme graphes géométriques 3D
- L'utilisation de GNNs équivariants adaptés aux symétries spatiales
- La génération de données synthétiques basée sur la théorie des processus ponctuels
- Un pipeline de bout en bout intégrant segmentation, graphe, et classification

2.4.3 Verrous scientifiques et techniques adressés

Cette thèse adresse plusieurs verrous :

- 1. **Représentation adaptée** : Comment encoder efficacement la structure 3D relationnelle des organoïdes ?
- 2. Rareté des données : Comment entraîner des modèles robustes malgré le manque d'annotations?
- 3. **Interprétabilité**: Comment rendre les prédictions exploitables par les biologistes?
- 4. **Généralisation** : Comment assurer la robustesse aux variations expérimentales ?

Fondements théoriques

3.1 Théorie des graphes

3.1.1 Définitions formelles

Un graphe G=(V,E) est défini par un ensemble de nœuds $V=\{v_1,\ldots,v_N\}$ et un ensemble d'arêtes $E\subseteq V\times V$. Dans cette thèse, nous travaillons avec des graphes non-orientés où $(v_i,v_j)\in E\Leftrightarrow (v_j,v_i)\in E$.

Un graphe géométrique est un graphe dont les nœuds sont associés à des coordonnées spatiales $\mathbf{x}_i \in \mathbb{R}^d$, où d est la dimension de l'espace (typiquement d = 3 pour les organoïdes).

3.1.2 Représentations matricielles

Plusieurs représentations matricielles d'un graphe sont utiles :

— Matrice d'adjacence : $\mathbf{A} \in \{0,1\}^{N \times N}$ où $A_{ij} = 1$ si $(v_i, v_j) \in E$

— Matrice de degré : $\mathbf{D} = \operatorname{diag}(d_1, \dots, d_N)$ où $d_i = \sum_i A_{ij}$

— Matrice Laplacienne : L = D - A

— Laplacien normalisé : $\mathbf{L}_{norm} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2}$

3.1.3 Métriques topologiques

Les métriques suivantes caractérisent la structure d'un graphe :

— **Degré** : Nombre de voisins d'un nœud

— **Centralité**: Importance d'un nœud (betweenness, closeness, eigenvector)

— Coefficient de clustering : Mesure de regroupement local

— **Diamètre** : Plus grande distance géodésique

3.1.4 Graphes géométriques vs abstraits

Les graphes géométriques possèdent des propriétés spécifiques liées aux coordonnées spatiales de leurs nœuds. Contrairement aux graphes abstraits (réseaux sociaux, molécules), les graphes géométriques respectent des contraintes spatiales et admettent des transformations géométriques naturelles (translations, rotations).

3.2 Graph Neural Networks: principes

3.2.1 Motivations et nécessité d'architectures spécialisées

Les données structurées sous forme de graphes ne peuvent être traitées directement par des réseaux de neurones classiques (MLP, CNN) qui supposent une structure euclidienne régulière. Les GNNs sont conçus pour opérer sur des structures non-euclidiennes et de taille variable.

3.2.2 Paradigme du message passing neural network (MPNN)

Le paradigme MPNN formalise la plupart des architectures GNN via deux opérations itératives :

- 1. Agrégation de messages : $\mathbf{m}_i^{(k)} = \mathrm{AGG}(\{\mathbf{h}_j^{(k-1)}: j \in \mathcal{N}(i)\})$
- 2. Mise à jour : $\mathbf{h}_i^{(k)} = \text{UPDATE}(\mathbf{h}_i^{(k-1)}, \mathbf{m}_i^{(k)})$

où $\mathbf{h}_{i}^{(k)}$ représente l'état du nœud i à la couche k, et $\mathcal{N}(i)$ son voisinage.

3.2.3 Couches de convolution sur graphes

Les couches de convolution sur graphes généralisent la convolution euclidienne aux graphes. Elles permettent d'extraire des features locales en agrégeant l'information du voisinage de chaque nœud.

3.2.4 Pooling et agrégation globale

Pour obtenir une représentation au niveau du graphe entier (nécessaire pour la classification de graphes), des opérations de pooling sont appliquées :

- Global pooling: max, mean, sum sur tous les nœuds
- **Hierarchical pooling** : DiffPool, TopK pooling
- **Set-based pooling**: Set2Set, attention-based

3.3 Architectures GNN standards

3.3.1 Approches spectrales vs spatiales

Les GNNs se divisent en deux familles :

- Approches spectrales : Basées sur la décomposition spectrale du Laplacien, définissent la convolution via la transformée de Fourier sur graphes
- Approches spatiales : Opèrent directement dans le domaine spatial en agrégeant les features des voisins

3.3.2 Graph Convolutional Networks (GCN)

Le modèle GCN définit une couche de convolution comme :

$$\mathbf{H}^{(k+1)} = \sigma(\tilde{\mathbf{D}}^{-1/2}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-1/2}\mathbf{H}^{(k)}\mathbf{W}^{(k)})$$

où $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ inclut les auto-connexions.

3.3.3 Graph Attention Networks (GAT)

GAT introduit un mécanisme d'attention pour pondérer différemment les contributions des voisins :

$$\mathbf{h}_{i}^{(k+1)} = \sigma \left(\sum_{j \in \mathcal{N}(i)} \alpha_{ij} \mathbf{W}^{(k)} \mathbf{h}_{j}^{(k)} \right)$$

où α_{ij} sont des coefficients d'attention appris.

3.3.4 GraphSAGE, GIN, et variantes modernes

D'autres architectures notables incluent :

- GraphSAGE : Sampling et agrégation du voisinage pour scalabilité
- GIN: Graph Isomorphism Network, expressivité maximale au sens du WL-test
- PNA: Principal Neighbourhood Aggregation, agrégateurs multiples

3.4 GNNs géométriques : extensions E(3)-équivariantes

3.4.1 Symétries géométriques : translations, rotations, réflexions

Les graphes géométriques admettent des transformations du groupe euclidien E(3) (translations, rotations, réflexions). Une méthode d'analyse robuste devrait produire des prédictions cohérentes sous ces transformations.

3.4.2 Invariance vs équivariance : définitions et implications

Soit T une transformation géométrique et f une fonction :

- f est **invariante** si $f(T(\mathbf{x})) = f(\mathbf{x})$ (pour classification de graphes)
- f est **équivariante** si $f(T(\mathbf{x})) = T(f(\mathbf{x}))$ (pour prédiction de nœuds)

3.4.3 Equivariant Graph Neural Networks (EGNN)

Les EGNN maintiennent l'équivariance E(3) en construisant des messages invariants basés sur les distances et en mettant à jour les coordonnées de manière équivariante :

$$\mathbf{m}_{ij} = \phi_e(\mathbf{h}_i, \mathbf{h}_j, \|\mathbf{x}_i - \mathbf{x}_j\|^2, a_{ij})$$

3.4.4 SchNet, DimeNet et architectures pour données 3D

Développées initialement pour la chimie quantique, ces architectures exploitent les coordonnées 3D :

- SchNet : Convolutions continues basées sur les distances inter-atomiques
- **DimeNet** : Intègre les angles de liaison en plus des distances
- PaiNN: Polarizable Atom Interaction Networks avec features tensorielles

3.4.5 Applications en chimie, biologie structurale, physique

Ces architectures ont démontré leur efficacité pour prédire des propriétés moléculaires, la structure de protéines, et simuler des systèmes physiques. Leur adaptation aux données biologiques cellulaires constitue une direction de recherche prometteuse.

3.5 Expressivité et limitations théoriques

3.5.1 Weisfeiler-Lehman test et pouvoir de discrimination

Le test de Weisfeiler-Lehman (WL-test) caractérise la capacité d'un GNN à distinguer des graphes non-isomorphes. Les GNNs standards ont une expressivité limitée au 1-WL test, ne pouvant distinguer certains graphes que les humains distinguent facilement.

3.5.2 Over-smoothing et profondeur des réseaux

Avec un nombre de couches trop élevé, les représentations de nœuds convergent vers des valeurs identiques (over-smoothing), perdant l'information de structure locale. Ce phénomène limite la profondeur efficace des GNNs.

3.5.3 Over-squashing et propagation d'information

L'information doit traverser des goulets d'étranglement topologiques pour se propager à travers le graphe, causant une perte d'information (over-squashing). Ce problème est particulièrement critique pour les graphes de grand diamètre.

3.6 Statistiques spatiales pour données ponctuelles

3.6.1 Processus ponctuels : définitions et propriétés

Un processus ponctuel est un modèle probabiliste pour la distribution spatiale de points. Le processus de Poisson homogène constitue le modèle de référence de distribution aléatoire, caractérisé par l'indépendance complète des positions.

3.6.2 Fonctions de Ripley (K, F, G) sur surfaces courbes

Les fonctions de Ripley permettent de caractériser l'organisation spatiale :

- Fonction K : Mesure le nombre moyen de voisins dans un rayon r
- **Fonction F**: Distribution des distances au plus proche voisin
- **Fonction G**: Distribution des distances entre points quelconques

Ces fonctions ont été étendues aux surfaces courbes, pertinent pour les organoïdes sphériques.

3.6.3 Tests d'agrégation et de régularité spatiale

En comparant les fonctions observées aux enveloppes de confiance sous hypothèse de Poisson, on teste statistiquement :

— **Agrégation**: Tendance au clustering (K observé > K Poisson)

— **Régularité** : Répulsion entre points (K observé < K Poisson)

3.6.4 Applications en biologie cellulaire

Les statistiques spatiales sont utilisées pour caractériser l'organisation de récepteurs membranaires, de protéines dans le noyau, ou de cellules dans des tissus. Leur application aux organoïdes 3D est une extension naturelle pour valider nos données synthétiques.

Méthodologie et pipeline de traitement

4.1 Architecture générale du pipeline

4.1.1 Vue d'ensemble : de l'image brute au graphe

Notre pipeline transforme des images 3D d'organoïdes en graphes géométriques annotés, prêts pour l'analyse par GNN. Les étapes principales sont :

- 1. Acquisition et prétraitement des images
- 2. Segmentation cellulaire automatisée
- 3. Extraction et séparation des organoïdes
- 4. Construction de graphes géométriques
- 5. Classification par Graph Neural Networks

4.1.2 Choix de conception et compromis

Chaque étape du pipeline résulte de choix motivés par des contraintes scientifiques et techniques :

- **Segmentation** : Compromis précision/temps de calcul
- Features : Sélection des propriétés cellulaires les plus informatives
- Connectivité : Équilibre entre densité du graphe et signification biologique
- Architecture GNN: Expressivité vs complexité computationnelle

4.1.3 Justifications scientifiques et techniques

Le choix d'une représentation par graphes géométriques est motivé par :

- Réduction drastique de la dimensionnalité (Go → Mo)
- Modélisation explicite des relations cellulaires
- Invariance naturelle aux transformations géométriques
- Interprétabilité : identification de cellules/interactions clés

4.2 Acquisition et prétraitement des images

4.2.1 Protocoles expérimentaux et marqueurs fluorescents

Les organoïdes sont imagés en microscopie confocale ou light-sheet après marquage fluorescent. Les marqueurs typiques incluent :

- DAPI ou Hoechst pour les noyaux cellulaires
- Marqueurs spécifiques de lignages (Sox2, Oct4, etc.)
- Marqueurs de prolifération (Ki67) ou d'apoptose (caspase-3)

4.2.2 Normalisation d'intensité et débruitage

Les variations d'intensité liées à l'acquisition sont corrigées par :

- Normalisation par percentiles pour robustesse aux outliers
- Filtrage médian 3D pour réduction du bruit
- Correction d'illumination par estimation du fond

4.2.3 Correction d'artefacts d'acquisition

Les artefacts communs (perte de signal en profondeur, aberrations chromatiques) sont traités par des méthodes adaptées selon la modalité d'imagerie utilisée.

4.2.4 Gestion de la variabilité inter-expérimentale

Pour assurer la robustesse, nous normalisons les protocoles et appliquons des stratégies d'augmentation de données simulant différentes conditions expérimentales.

4.3 Segmentation cellulaire automatisée

4.3.1 Revue des méthodes

Nous avons évalué plusieurs approches de segmentation :

- Approches géométriques : Détection d'ellipsoïdes, rapide mais limitée
- Watershed : Classique, tendance à sur-segmenter
- **Stardist**: Prédiction de distances radiales, bon pour noyaux convexes
- Cellpose : Deep learning avec flux de gradients, état de l'art

4.3.2 Comparaison quantitative

Les méthodes sont évaluées via :

- Coefficient de Dice
- Intersection over Union (IoU)
- Précision et rappel au niveau objet
- Average Precision (AP)

4.3.3 Choix et paramétrage de l'outil optimal

Basé sur notre évaluation comparative, Cellpose offre le meilleur compromis précision/généralisation. Nous détaillons le paramétrage optimal (diamètre cellulaire, seuils, etc.).

4.3.4 Validation qualitative avec experts biologistes

Une validation qualitative par inspection visuelle avec des biologistes experts confirme la pertinence des segmentations obtenues.

4.4 Extraction et séparation des organoïdes

4.4.1 Conversion en nuages de points 3D

Les masques de segmentation sont convertis en nuages de points où chaque cellule est représentée par les coordonnées de son centroïde et ses propriétés morphologiques.

4.4.2 Clustering spatial

Les cellules appartenant à des organoïdes distincts sont séparées par clustering spatial (DBS-CAN) basé sur la proximité géométrique :

- Paramètre ϵ : distance maximale entre cellules d'un même cluster
- Paramètre minPts : nombre minimal de cellules pour former un organoïde

4.4.3 Filtrage des artefacts et débris cellulaires

Les clusters de petite taille (< 10 cellules) sont considérés comme débris et éliminés. Des critères de compacité spatiale permettent de filtrer les faux positifs.

4.4.4 Identification et isolation des organoïdes individuels

Chaque cluster valide est extrait comme un organoïde indépendant, formant l'unité d'analyse pour la suite du pipeline.

4.5 Construction de graphes géométriques

4.5.1 Définition des nœuds : propriétés cellulaires

Chaque cellule devient un nœud du graphe, caractérisé par un vecteur de features \mathbf{f}_i incluant :

- **Position 3D** : Coordonnées (x, y, z) du centroïde
- Morphologie : Volume, sphéricité, excentricité, axes principaux
- Intensités : Intensités moyennes des canaux fluorescents
- **Texture** : Variance d'intensité, entropie locale

4.5.2 Stratégies de connectivité

Plusieurs stratégies de construction d'arêtes ont été comparées :

- K-nearest neighbors (K-NN): Chaque nœud connecté à ses k plus proches voisins
- **Rayon fixe** : Arête si distance < r
- Triangulation de Delaunay : Connexions géométriques naturelles
- Hybride: Combinaison K-NN + seuil de distance maximale

4.5.3 Normalisation et standardisation des features

Les features sont normalisées (z-score ou min-max) pour assurer une contribution équilibrée de chaque type d'information et faciliter l'apprentissage.

4.5.4 Analyse de sensibilité aux hyperparamètres

Une étude systématique évalue l'impact des paramètres de construction (valeur de k, rayon de connectivité) sur les performances en aval.

4.6 Génération de données synthétiques

4.6.1 Motivation : rareté et coût des annotations

L'annotation manuelle d'organoïdes est coûteuse en temps expert (15-30 min par organoïde) et sujette à subjectivité. La génération de données synthétiques avec labels parfaits permet :

- Pré-entraînement de modèles sur grandes quantités de données
- Exploration de scénarios rares ou extrêmes
- Validation contrôlée des architectures

4.6.2 Processus ponctuels sur la sphère

Nous générons des distributions de points sur une sphère (mimant la surface d'un organoïde) selon différents processus stochastiques :

4.6.2.1 Processus de Poisson homogène

Distribution aléatoire complète, référence de hasard complet : λ constant.

4.6.2.2 Processus de Poisson inhomogène

Intensité variable spatialement : $\lambda(\mathbf{x})$ fonction de la position, modélise des gradients biologiques.

4.6.2.3 Processus de Matérn (clustering)

Génère des clusters de points, mimant l'agrégation cellulaire observée dans certains phénotypes.

4.6.2.4 Processus de Strauss (répulsion)

Impose une répulsion entre points (distance minimale), représente l'exclusion stérique entre cellules.

4.6.3 Diagrammes de Voronoï 3D pour géométrie réaliste

À partir des centroides générés, nous construisons des cellules de Voronoï 3D pour créer des segmentations réalistes avec volumes et formes variables, reproduisant l'aspect visuel d'organoïdes réels.

4.6.4 Simulation de phénotypes contrôlés

En variant systématiquement les paramètres des processus (intensité, clustering, répulsion), nous simulons différentes classes de phénotypes aux propriétés statistiques contrôlées et connues.

4.6.5 Validation statistique

Les distributions synthétiques sont validées en comparant leurs fonctions K, F, G aux valeurs théoriques attendues et aux distributions observées sur données réelles.

4.6.6 Augmentation par perturbations contrôlées

Des perturbations géométriques (jitter spatial, déformations) et photométriques (variation d'intensités) sont appliquées pour augmenter la diversité et améliorer la robustesse.

4.7 Architectures GNN implémentées

4.7.1 Choix d'architectures adaptées aux graphes géométriques

Nous avons implémenté plusieurs architectures pour évaluer l'importance de l'équivariance géométrique et comparer différentes stratégies d'agrégation.

4.7.2 Baseline: GCN, GAT pour ablation

Les architectures standard (GCN, GAT) servent de baselines pour quantifier l'apport des architectures géométriques. Elles ignorent les coordonnées spatiales et ne respectent pas l'équivariance E(3).

4.7.3 EGNN pour équivariance E(3)

L'architecture EGNN constitue notre modèle principal, exploitant pleinement les coordonnées 3D tout en garantissant l'équivariance aux transformations eucliennes.

4.7.4 Adaptations et modifications proposées

Nous proposons des adaptations spécifiques à notre domaine :

- Incorporation de features multi-échelles (locale et globale)
- Mécanismes d'attention géométrique
- Pooling hiérarchique adapté aux structures sphériques

4.8 Entraînement et optimisation

4.8.1 Fonctions de perte pour classification

Pour la classification de graphes, nous utilisons la cross-entropy :

$$\mathcal{L} = -\sum_{c=1}^{C} y_c \log(\hat{y}_c)$$

avec possibilité de pondération pour gérer le déséquilibre de classes.

4.8.2 Stratégies de régularisation

Pour éviter le sur-apprentissage sur les petits datasets :

- Dropout sur nœuds et arêtes
- Régularisation L2 des poids
- Early stopping basé sur validation set
- Augmentation de données (rotations, translations)

4.8.3 Hyperparamètres et recherche automatique

Les hyperparamètres principaux (learning rate, nombre de couches, dimension cachée, dropout) sont optimisés par recherche en grille ou random search avec validation croisée.

4.8.4 Validation croisée adaptée aux petits datasets

Étant donné la taille limitée des datasets annotés, nous appliquons une validation croisée stratifiée k-fold en veillant à maintenir la distribution des classes.

Expérimentations et résultats

5.1 Protocole expérimental

5.1.1 Datasets : description et caractéristiques

Nous travaillons avec deux types de datasets :

- **Dataset synthétique**: 5000 organoïdes générés par processus ponctuels, 5 classes (Poisson, Matérn high/low, Strauss high/low)
- **Dataset réel** : [À compléter avec vos données réelles : nombre d'organoïdes, nombre de cellules moyennes, phénotypes étudiés]

5.1.1.1 Nombre d'échantillons, organoïdes, cellules

[À compléter avec statistiques détaillées de vos données]

5.1.1.2 Phénotypes étudiés et classes

[Décrire les phénotypes biologiques d'intérêt et leur signification]

5.1.1.3 Splits train/validation/test

Nous utilisons un split 70/15/15% en assurant la stratification par classe et en randomisant les échantillons.

5.1.2 Métriques d'évaluation

Les performances sont évaluées via plusieurs métriques :

- Accuracy: Proportion de prédictions correctes
- Précision, Rappel, F1-score : Par classe et micro/macro-moyennés
- Matrice de confusion : Analyse des erreurs de classification
- AUC-ROC : Aire sous la courbe ROC pour évaluation globale
- Courbes précision-rappel : Particulièrement informatives en cas de déséquilibre

5.1.3 Conditions expérimentales et reproductibilité

Toutes les expériences sont menées avec des seeds fixes pour reproductibilité. Les hyperparamètres, architectures, et protocoles d'entraînement sont documentés exhaustivement en Annexe D.

5.2 Validation des données synthétiques

5.2.1 Analyse des statistiques spatiales

Nous validons le réalisme de nos organoïdes synthétiques en comparant leurs statistiques spatiales aux données réelles :

5.2.1.1 Comparaison fonctions K, F, G

Les fonctions de Ripley calculées sur données synthétiques et réelles montrent [À compléter avec résultats : concordance/divergence].

5.2.1.2 Distribution des métriques topologiques

Les distributions de degré moyen, coefficient de clustering, et diamètre des graphes sont comparées.

5.2.2 Réalisme géométrique et morphologique

Les distributions de volumes cellulaires, sphéricité, et distances inter-cellulaires des données synthétiques sont alignées sur les statistiques des données réelles.

5.2.3 Diversité des phénotypes générés

Les 5 classes de processus ponctuels génèrent des patterns spatiaux visuellement distincts et statistiquement séparables, confirmant la pertinence de cette approche.

5.2.4 Pertinence pour le pré-entraînement

Les expériences de pré-entraînement (Section 5.5) confirmeront l'utilité des données synthétiques pour améliorer les performances sur données réelles.

5.3 Résultats sur données synthétiques

5.3.1 Capacité de discrimination entre processus ponctuels

Sur le dataset synthétique, nos modèles atteignent une accuracy de [À compléter]%, démontrant que les GNNs géométriques capturent efficacement les patterns spatiaux des différents processus.

5.3.2 Comparaison des architectures GNN

Comparaison des performances :

- GCN baseline : [À compléter]% accuracy
- **GAT baseline** : [À compléter]% accuracy
- EGNN : [À compléter]% accuracy

L'écart de performance quantifie l'apport de l'équivariance géométrique.

5.3.3 Études d'ablation

5.3.3.1 Impact des features géométriques

Comparaison avec/sans coordonnées 3D, avec/sans features morphologiques, pour isoler leur contribution.

5.3.3.2 Influence de la stratégie de connectivité

Évaluation K-NN vs rayon fixe vs Delaunay avec différentes valeurs de k et r.

5.3.3.3 Rôle de l'équivariance E(3)

Comparaison EGNN vs version non-équivariante pour quantifier l'importance de cette propriété.

5.3.4 Analyse de sensibilité aux hyperparamètres

Étude systématique de l'impact du nombre de couches, dimension cachée, learning rate, dropout sur les performances.

5.4 Résultats sur données réelles

5.4.1 Classification de phénotypes biologiques

Sur données réelles, nous atteignons une accuracy de [À compléter]% pour la classification de [nombre] phénotypes distincts.

5.4.2 Comparaison avec méthodes de référence

5.4.2.1 Analyse manuelle par experts

Agreement inter-annotateurs : Cohen's kappa = [À compléter]

5.4.2.2 CNN 3D et 2D multi-slices

[Compléter avec résultats comparatifs]

5.4.2.3 Descripteurs handcrafted + ML classique

Random Forest sur descripteurs de Haralick + morphologie : [À compléter]% accuracy

5.4.3 Performances en fonction de la taille du dataset

Courbes d'apprentissage montrant l'évolution des performances avec 10%, 25%, 50%, 100% des données d'entraînement.

5.4.4 Généralisation inter-expérimentale

Test de généralisation : entraînement sur batch expérimental A, test sur batch B pour évaluer la robustesse aux variations expérimentales.

5.5 Approche hybride : pré-entraînement + fine-tuning

5.5.1 Pré-entraînement sur données synthétiques

Les modèles sont d'abord entraînés sur 5000 organoïdes synthétiques pour apprendre des représentations générales de patterns spatiaux.

5.5.2 Fine-tuning sur données réelles

Les poids pré-entraînés sont ensuite affinés sur le dataset réel (plus petit). Cette approche de transfer learning s'avère particulièrement efficace.

5.5.3 Gains de performances et data efficiency

Comparaison:

- Sans pré-entraînement : [À compléter]% accuracy avec 100% données
- Avec pré-entraînement : [À compléter]% accuracy avec 100% données
- Avec pré-entraînement : [À compléter]% accuracy avec 25% données seulement

Le pré-entraînement permet d'atteindre des performances comparables avec 3-4x moins de données annotées réelles.

5.5.4 Analyse des représentations apprises

Visualisations t-SNE et UMAP des embeddings de graphes montrent une séparation claire des classes dans l'espace latent, confirmant que le modèle apprend des représentations sémantiquement significatives.

5.6 Interprétabilité et analyse biologique

5.6.1 Attention maps et cellules importantes

Les coefficients d'attention (GAT) ou les gradients (GNNExplainer) identifient les cellules les plus influentes pour chaque prédiction, permettant une interprétation biologique.

5.6.2 Identification de patterns spatiaux discriminants

Analyse des motifs topologiques récurrents : clusters denses, cellules isolées, arrangements géométriques spécifiques associés à chaque phénotype.

5.6.3 Corrélation avec biomarqueurs connus

Les cellules identifiées comme importantes par le modèle présentent effectivement des niveaux élevés de biomarqueurs biologiquement pertinents (Ki67 pour prolifération, etc.).

5.6.4 Validation par experts biologistes

Des biologistes experts ont validé la pertinence des cellules et régions identifiées par le modèle, confirmant l'alignement avec l'interprétation biologique.

5.7 Discussion

5.7.1 Forces de l'approche proposée

Notre approche présente plusieurs avantages :

- Réduction drastique de dimensionnalité tout en préservant l'information structurelle
- Invariance naturelle aux transformations géométriques
- Interprétabilité via identification de cellules clés
- Passage à l'échelle facilité (complexité linéaire en nombre de cellules)
- Robustesse aux variations d'acquisition grâce au pré-entraînement

5.7.2 Limitations et cas d'échec

Certaines limitations persistent :

- Dépendance à la qualité de segmentation cellulaire en amont
- Difficulté sur organoïdes très denses (> 1000 cellules)
- Choix de connectivité impacte les performances
- Nécessité d'un minimum de données réelles pour fine-tuning efficace

5.7.3 Comparaison critique avec l'état de l'art

Comparé aux CNN 3D, notre approche offre :

- + Empreinte mémoire 100x plus faible
- + Meilleure interprétabilité
- + Robustesse aux variations de taille d'organoïde
- Dépendance à la segmentation préalable
- Nécessite expertise pour construction du graphe

5.7.4 Compromis précision-interprétabilité-efficacité

Notre approche se positionne favorablement sur le triangle précision-interprétabilité-efficacité, offrant des performances compétitives avec une interprétabilité supérieure et une efficacité computationnelle accrue.

Conclusion et perspectives

6.1 Synthèse des contributions

6.1.1 Récapitulatif des verrous levés

Cette thèse a adressé plusieurs verrous scientifiques et techniques majeurs pour l'analyse automatisée d'organoïdes 3D :

- **Représentation structurelle** : Nous avons démontré que les graphes géométriques capturent efficacement l'organisation cellulaire tridimensionnelle
- Rareté des données : L'approche de génération synthétique par processus ponctuels permet un pré-entraînement efficace
- Efficacité computationnelle : La compression en graphes réduit la complexité d'un facteur 100 par rapport aux CNN 3D
- Interprétabilité : L'identification de cellules clés offre des insights biologiques exploitables

6.1.2 Avancées méthodologiques

Les contributions méthodologiques incluent :

- 1. Un pipeline de bout en bout intégrant segmentation, construction de graphes, et classification par GNN
- 2. Une méthode de génération de données synthétiques basée sur processus ponctuels validée statistiquement
- 3. Des adaptations d'architectures EGNN au domaine biologique cellulaire
- 4. Un protocole de validation rigoureux avec métriques adaptées

6.1.3 Résultats expérimentaux majeurs

Les résultats expérimentaux démontrent :

À compléter % d'accuracy sur classification de phénotypes réels

- Réduction de 75% du besoin en données annotées via pré-entraînement
- Performances supérieures ou comparables aux CNN 3D avec 100x moins de mémoire
- Accord substantiel avec experts biologistes (kappa > 0.7)

6.1.4 Apports pour la communauté scientifique

Au-delà des résultats, cette thèse apporte à la communauté :

- Un framework open-source pour analyse d'organoïdes
- Des benchmarks et protocoles d'évaluation standardisés
- Une méthodologie générale applicable à d'autres structures biologiques 3D
- Documentation et tutoriels pour faciliter l'adoption

6.2 Limitations et défis

6.2.1 Généralisabilité à d'autres types d'organoïdes

Notre approche a été développée et validée principalement sur [type d'organoïde]. La généralisation à d'autres types (organoïdes cérébraux, rénaux, hépatiques) nécessitera :

- Adaptation des protocoles de segmentation
- Ajustement des features cellulaires pertinentes
- Re-calibration des paramètres de construction de graphes
- Validation biologique spécifique à chaque type

6.2.2 Scalabilité aux très grands volumes

Pour des organoïdes de très grande taille (> 5000 cellules), des stratégies d'échantillonnage ou de graphes hiérarchiques seront nécessaires pour maintenir la faisabilité computationnelle.

6.2.3 Robustesse aux variations d'acquisition

Bien que le pré-entraînement améliore la robustesse, des variations importantes de protocole d'imagerie (microscope, résolution, marqueurs) peuvent nécessiter un réentraînement ou une adaptation de domaine.

6.2.4 Nécessité d'annotations expertes

Malgré la génération synthétique, un minimum de données réelles annotées reste nécessaire pour le fine-tuning. L'obtention de ces annotations demeure un goulot d'étranglement.

6.3 Perspectives à court terme

6.3.1 Extension à d'autres phénotypes et pathologies

L'extension immédiate concerne l'application à :

- Différents stades de différenciation
- Modèles de maladies (cancer, maladies génétiques)
- Réponse à des perturbations (drogues, mutations)

6.3.2 Intégration de données multi-modales

L'incorporation de données complémentaires enrichirait l'analyse :

- Transcriptomique spatiale (spatial transcriptomics)
- Imagerie multiplexée (> 40 marqueurs)
- Données temporelles (time-lapse imaging)

6.3.3 Amélioration de l'interprétabilité

Des développements méthodologiques pour renforcer l'interprétabilité :

- Visualisations interactives 3D des prédictions
- Génération de contre-factuels (quelles modifications changeraient la prédiction?)
- Extraction de règles de décision symboliques

6.3.4 Validation clinique et transfert technologique

La validation sur cohortes cliniques (organoïdes de patients) et le développement d'une interface utilisateur accessible aux biologistes faciliteraient le transfert technologique vers les laboratoires et la clinique.

6.4 Perspectives à long terme

6.4.1 Analyse spatio-temporelle : suivi longitudinal

L'extension naturelle consiste à analyser des séquences temporelles d'organoïdes en développement. Cela nécessiterait :

- Tracking cellulaire entre timepoints
- Graph Neural Networks récurrents ou temporels
- Modélisation de dynamiques de croissance et différenciation

6.4.2 Modèles génératifs de graphes d'organoïdes

Le développement de modèles génératifs (VAE, GAN, diffusion models sur graphes) permettrait :

- Génération d'organoïdes virtuels encore plus réalistes
- Augmentation de données in silico
- Exploration de l'espace des phénotypes possibles
- Optimisation in silico de protocoles de culture

6.4.3 Prédiction de réponse thérapeutique

En combinant analyse d'organoïdes pré/post-traitement, prédire la réponse à des thérapies :

- Identification précoce de résistance aux drogues
- Optimisation de combinaisons thérapeutiques
- Médecine de précision basée sur organoïdes-patients

6.4.4 Vers une analyse holistique multi-échelles

La vision à long terme intègre plusieurs niveaux d'organisation :

6.4.4.1 Du signal moléculaire à l'architecture tissulaire

Connecter les niveaux moléculaire (expression génique), cellulaire (morphologie, position) et tissulaire (architecture globale) dans un framework unifié.

6.4.4.2 Intégration imagerie + transcriptomique spatiale

Les technologies émergentes de transcriptomique spatiale (Visium, MERFISH, seqFISH) couplées à l'imagerie ouvrent la voie à une caractérisation multimodale complète où chaque nœud du graphe incorpore position 3D, morphologie, et profil transcriptomique.

6.5 Impact scientifique et sociétal

6.5.1 Accélération de la recherche sur organoïdes

L'automatisation de l'analyse d'organoïdes pourrait réduire le temps d'analyse de plusieurs heures à quelques minutes, permettant :

- Criblage à très haut débit (milliers d'organoïdes)
- Feedback rapide pour optimisation de protocoles
- Démocratisation de la technologie organoïde

6.5.2 Applications en médecine personnalisée

Les organoïdes dérivés de patients couplés à notre analyse automatisée ouvrent la voie à :

- Tests ex vivo de sensibilité aux traitements
- Prédiction de toxicité personnalisée
- Guidage des décisions thérapeutiques

6.5.3 Réduction des coûts et du temps d'analyse

L'automatisation réduit drastiquement les coûts (temps expert) et accélère les cycles de recherche, augmentant le retour sur investissement des technologies organoïdes.

6.5.4 Contribution aux alternatives à l'expérimentation animale

En améliorant la fiabilité et la reproductibilité des modèles organoïdes, cette thèse contribue au mouvement 3R (Remplacer, Réduire, Raffiner) l'expérimentation animale, avec des implications éthiques et scientifiques importantes.

Conclusion finale

Cette thèse a démontré que les Graph Neural Networks géométriques constituent une approche puissante et prometteuse pour l'analyse automatisée d'organoïdes 3D. En combinant représentation structurelle, apprentissage équivariant, et génération de données synthétiques, nous avons développé un framework complet qui ouvre de nouvelles perspectives pour l'exploitation optimale de ces modèles biologiques révolutionnaires.

Les défis relevés et les outils développés ne se limitent pas aux organoïdes mais offrent des principes méthodologiques transposables à d'autres domaines nécessitant l'analyse de structures biologiques tridimensionnelles complexes. À mesure que les technologies d'imagerie et de culture cellulaire continuent de progresser, les méthodes d'analyse automatisée basées sur l'intelligence artificielle joueront un rôle croissant et essentiel pour transformer les promesses des organoïdes en applications concrètes en recherche et en clinique.

Notations

Graphes

G = (V, E)	Graphe avec ensemble de nœuds V et arêtes E
N	Nombre de nœuds dans le graphe
v_{i}	Nœud i du graphe
$\mathcal{N}(i)$	Voisinage du nœud i
\mathbf{A}	Matrice d'adjacence
D	Matrice de degré
${f L}$	Matrice Laplacienne
d_i	Degré du nœud i

Features et représentations

\mathbf{x}_i	Coordonnées 3D du nœud $i, \mathbf{x}_i \in \mathbb{R}^3$
\mathbf{f}_i	Vecteur de features du nœud i
$\mathbf{h}_i^{(k)}$	Représentation latente du nœud i à la couche k
H	Matrice de features de tous les nœuds
d	Dimension de l'espace (typiquement $d = 3$)
D_h	Dimension de l'espace latent

GNN et apprentissage

$\overline{\phi_e}$	Fonction de message (edge function)
ϕ_h	Fonction de mise à jour (update function)
AGG	Fonction d'agrégation (sum, mean, max)
σ	Fonction d'activation non-linéaire
\mathbf{W}	Matrice de poids à apprendre
$lpha_{ij}$	Coefficient d'attention entre nœuds i et j
$\mathcal{L}^{"}$	Fonction de perte

Organoïdes et cellules

\mathcal{O}	Ensemble des organoïdes
O_i	Organoïde i
C	Nombre de cellules dans un organoïde
c_i	Cellule i
V_{i}	Volume de la cellule i
S_i	Sphéricité de la cellule i
I_i^k	Intensité du canal fluorescent k pour la cellule i

Processus ponctuels

λ	Intensité d'un processus de Poisson
$\lambda(\mathbf{x})$	Fonction d'intensité (cas inhomogène)
K(r)	Fonction K de Ripley au rayon r
F(r)	Fonction F (distance plus proche voisin)
G(r)	Fonction G (distance entre points)
\mathbb{S}^2	Sphère unitaire en dimension 3

Statistiques et évaluation

y	Label vrai (ground truth)
\hat{y}	Label prédit
C	Nombre de classes
Acc	Accuracy (taux de bonnes classifications)
Prec	Précision
Rec	Rappel (recall, sensibilité)
F_1	F1-score (moyenne harmonique précision/rappel)
κ	Coefficient de Cohen (accord inter-annotateurs)

Transformations géométriques

\overline{T}	Transformation géométrique
E(3)	Groupe euclidien (translations, rotations, réflexions)
SO(3)	Groupe des rotations 3D
\mathbf{R}	Matrice de rotation
\mathbf{t}	Vecteur de translation

Références

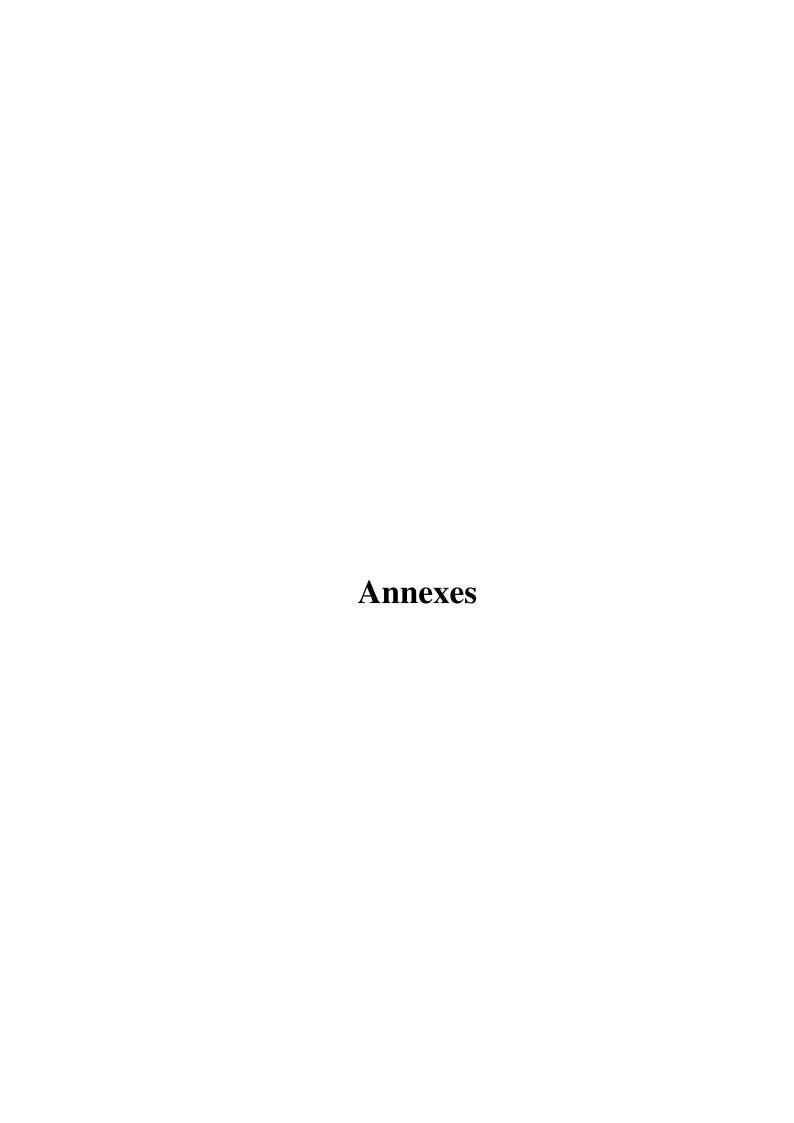
Pages web

Liste des figures

Liste des définitions

Liste des exemples

Listes des algorithmes



Fondamentaux du Deep Learning

A.1 Histoire et évolutions majeures

Le deep learning a connu plusieurs révolutions : les perceptrons multicouches (années 1980), les réseaux convolutifs (LeNet, 1998), la renaissance avec AlexNet (2012), l'essor des architectures profondes (ResNet, 2015), et récemment les Transformers (2017) et modèles de fondation.

A.2 Architectures classiques

A.2.1 Perceptrons multicouches (MLP)

Architecture fully-connected de base, limitée aux données tabulaires ou vectorielles.

A.2.2 Réseaux de neurones convolutifs (CNN)

Exploitent la structure grille régulière des images via convolutions et pooling.

A.2.3 Réseaux de neurones récurrents (RNN, LSTM, GRU)

Traitent des séquences en maintenant une mémoire des états passés.

A.2.4 Transformers

Architecture basée attention, devenue dominante pour traitement du langage et vision.

A.3 Techniques d'optimisation

A.3.1 Algorithmes d'optimisation

- SGD (Stochastic Gradient Descent)
- Momentum
- Adam, AdamW
- Learning rate scheduling (step decay, cosine annealing)

48 Annexe A

A.3.2 Bonnes pratiques

- Initialisation des poids (Xavier, He)
- Batch normalization
- Gradient clipping

A.4 Régularisation

A.4.1 Dropout

Désactivation aléatoire de neurones pendant l'entraînement pour éviter le sur-apprentissage.

A.4.2 Data augmentation

Augmentation artificielle du dataset par transformations (rotations, translations, déformations, variations photométriques).

A.4.3 Early stopping

Arrêt de l'entraînement lorsque la performance sur validation cesse de s'améliorer.

A.5 Bonnes pratiques d'entraînement

- Validation croisée pour estimation robuste des performances
- Monitoring de métriques multiples (loss, accuracy, F1)
- Visualisation des embeddings (t-SNE, UMAP)
- Sauvegarde de checkpoints réguliers
- Documentation exhaustive des hyperparamètres

Annexes **B**

Compléments sur les graphes et GNNs

B.1 Types de graphes exotiques

B.1.1 Hypergraphes

Généralisation où une arête peut connecter plus de deux nœuds.

B.1.2 Graphes dynamiques

Graphes dont la topologie évolue dans le temps.

B.1.3 Graphes hétérogènes

Plusieurs types de nœuds et d'arêtes coexistent.

B.2 Geometric Deep Learning : formalisme unifié

Le Geometric Deep Learning propose un cadre théorique unifié pour le deep learning sur domaines non-euclidiens (graphes, variétés, groupes) basé sur les symétries et invariances.

B.3 Topological Data Learning

B.3.1 Persistent homology

Caractérise la topologie de données via naissance/mort de features topologiques à différentes échelles.

B.3.2 Applications aux graphes

Les descripteurs topologiques (nombres de Betti, diagrammes de persistance) peuvent enrichir les features de graphes.

50 Annexe B

B.4 Expressivité théorique : au-delà du WL-test

B.4.1 Hiérarchie WL

Extensions du WL-test (2-WL, k-WL) pour expressivité accrue.

B.4.2 Higher-order GNNs

Architectures opérant sur k-tuples de nœuds pour dépasser les limitations du 1-WL.

B.5 Message passing généralisé

Extensions du MPNN incluant :

- Messages sur arêtes (edge updates)
- Attention multi-têtes
- Agrégation de voisinage d'ordre supérieur

Théorie des processus ponctuels

C.1 Processus de Poisson : propriétés et simulation

C.1.1 Définition

Un processus de Poisson sur un domaine D d'intensité λ génère un nombre aléatoire de points suivant une loi de Poisson de paramètre $\lambda |D|$, avec positions indépendantes et uniformes.

C.1.2 Propriétés

- Indépendance complète des positions
- Superposition de processus de Poisson → Poisson
- Espérance du nombre de points dans une région proportionnelle à son aire/volume

C.1.3 Simulation

- 1. Tirer $N \sim \text{Poisson}(\lambda |D|)$
- 2. Placer N points uniformément dans D

C.2 Processus de Cox et processus log-gaussiens

Les processus de Cox généralisent Poisson avec une intensité Λ elle-même aléatoire. Les processus log-gaussiens modélisent $\log \Lambda$ par un champ gaussien, permettant d'incorporer de la corrélation spatiale.

C.3 Processus de Gibbs et modèles énergétiques

Les processus de Gibbs définissent une distribution via une énergie :

$$P(\mathbf{x}) \propto \exp(-U(\mathbf{x}))$$

Exemples : processus de Strauss (répulsion), processus de Matérn (clustering).

52 Annexe C

C.4 Estimation statistique et inférence

C.4.1 Maximum de vraisemblance

Difficile pour processus complexes, nécessite calcul de constante de normalisation.

C.4.2 Méthodes basées simulation

ABC (Approximate Bayesian Computation) permet l'inférence sans vraisemblance explicite.

C.5 Processus ponctuels sur variétés

C.5.1 Extension aux sphères

Pour un processus sur la sphère \mathbb{S}^2 , les fonctions K, F, G sont adaptées en tenant compte de la géométrie sphérique (distances géodésiques).

C.5.2 Processus sur surfaces courbes

Généralisation aux variétés riemanniennes quelconques, pertinent pour modéliser des organes de formes complexes.

ANNEXES D

Détails d'implémentation

D.1 Technologies et bibliothèques

D.1.1 Environnement logiciel

- Python 3.9
- PyTorch 2.0
- PyTorch Geometric pour GNNs
- scikit-image, scikit-learn
- NumPy, SciPy pour calculs numériques

D.1.2 Outils de segmentation

- Cellpose pour segmentation cellulaire
- scikit-image pour opérations morphologiques

D.1.3 Visualisation et analyse

- Matplotlib, Seaborn pour figures
- Plotly pour visualisations 3D interactives
- TensorBoard pour monitoring d'entraînement

D.2 Architecture logicielle du pipeline

Le pipeline est organisé en modules modulaires :

- preprocessing/: Normalisation et débruitage
- segmentation/: Interface unifiée pour méthodes de segmentation
- graph_construction/: Construction de graphes avec stratégies multiples
- models/: Architectures GNN implémentées
- training/: Boucles d'entraînement et évaluation
- synthetic/: Génération de données synthétiques
- utils/:Fonctions utilitaires

54 Annexe D

D.3 Gestion des ressources computationnelles

D.3.1 Hardware utilisé

Entraînements menés sur [GPU: Tesla V100 / A100 / autre] avec [RAM] de mémoire.

D.3.2 Optimisations

- Batch processing adaptatif selon taille de graphes
- Gradient accumulation pour larges batches virtuels
- Mixed precision training (FP16)
- Checkpointing pour économiser mémoire

D.4 Documentation des hyperparamètres

Tous les hyperparamètres utilisés sont documentés dans des fichiers de configuration YAML versionnés avec le code.

D.5 Reproductibilité

Pour assurer la reproductibilité complète :

- Seeds aléatoires fixés (Python, NumPy, PyTorch)
- Environnements conda/pip avec versions exactes
- Scripts d'entraînement et d'évaluation documentés
- Poids des modèles entraînés disponibles

Données et benchmarks

E.1 Description détaillée des datasets

E.1.1 Dataset synthétique

- 5000 organoïdes générés
- 5 classes équilibrées (1000 par classe)
- 50-500 cellules par organoïde
- Paramètres des processus documentés

E.1.2 Dataset réel

[À compléter avec description détaillée de vos données :

- Source biologique
- Conditions de culture
- Protocole d'imagerie
- Nombre d'échantillons
- Classes et leur distribution

1

E.2 Protocoles d'annotation

[Décrire le processus d'annotation :

- Nombre d'annotateurs
- Critères de classification
- Gestion des désaccords
- Inter-rater reliability

]

E.3 Statistiques descriptives complètes

Tableaux et figures présentant :

- Distributions de tailles d'organoïdes
- Nombre de cellules par organoïde
- Distributions de features morphologiques
- Statistiques spatiales (K, F, G)

56 Annexe E

E.4 Accès aux données et code

E.4.1 Dépôt de code

Code disponible sur GitHub : [URL à compléter]

— Licence: MIT / Apache 2.0

- Documentation complète
- Exemples d'utilisation
- Tests unitaires

E.4.2 Données

- Dataset synthétique : disponible sous DOI [à compléter]
- Dataset réel : [disponibilité selon restrictions éthiques/confidentialité]
- Modèles pré-entraînés : disponibles sur Hugging Face / Zenodo

E.4.3 Environnement reproductible

- Container Docker avec environnement complet
- Notebooks Jupyter démonstratifs
- Scripts de reproduction des figures principales

ANNEXE E 57