



UNIVERSITÉ SORBONNE UNIVERSITÉ

ECOLE DOCTORALE INFORMATIQUE, TÉLÉCOMMUNICATIONS ET ELECTRONIQUE - ED130

INRIA DE PARIS / ÉQUIPE ALMANACH

THÈSE DE DOCTORAT

Discipline : Informatique

Présentée par

Rian TOUCHENT-SAAD

Pour obtenir le grade universitaire de

DOCTEUR de l'UNIVERSITÉ SORBONNE UNIVERSITÉ

Improving Clinical Information Extraction with Public-Data Pretraining

Présentée et soutenue publiquement le DATE devant le jury composé de :

Mathieu CONSTANT	ATILF (CNRS / Université de Lorraine)	Rapporteur
Aurélie NÉVÉOL	CNRS, LISN, Université Paris-Saclay	Rapportrice
Natalia GRABAR	CNRS, Université de Lille	Rapportrice
Pierre ZWEIGENBAUM	CNRS, LISN, Université Paris-Saclay	Examinateur
Thierry CHARNOIS	LIPN, Université Sorbonne Paris Nord	Examinateur
Laurent ROMARY	Inria - ALMANACh	Directeur
Éric DE LA CLERGERIE	Inria - ALMANACh	Co-encadrant

ABSTRACT

CONTENTS

1	INTRODUCTION	1
I	RELATED WORKS	3
2	LANGUAGE MODELS	5
2.1	Introduction	5
2.2	From Statistical to Neural Language Models	5
2.2.1	N-gram Models	5
2.2.2	Neural LMs	5
2.2.3	RNNs	6
2.3	Transformer	6
2.3.1	Self-Attention	6
2.3.2	Positional Encodings	6
2.3.3	Encoder vs Decoder	6
2.4	Pretraining Objectives	7
2.4.1	CLM (Causal Language Modeling)	7
2.4.2	MLM (Masked Language Modeling)	7
2.4.3	Objectifs hybrides (2024-2025)	7
2.5	Scaling Laws & LLMs	8
2.5.1	Scaling	8
2.5.2	Post-ChatGPT	8
2.6	Continual Pretraining	8
2.6.1	Fondements et terminologie	8
2.6.2	Catastrophic Forgetting	9
2.6.3	Solutions recentes (2024-2025)	9
2.6.4	From scratch vs continual PT	10
2.6.5	Encoders vs decoders: dynamiques différentes	10
2.6.6	Tentatives antérieures en français biomedical	10
2.7	Modèles pretrained biomedical	11
2.7.1	Encoders biomedical anglais	11
2.7.2	Modèles français	11
2.7.3	Decoders biomedical	11
2.8	Architectures modernes	12
2.8.1	Modernisation des encoders	12
2.8.2	Efficient Attention	12
2.8.3	Long context	13

Contents

2.9	Tokenization	13
2.9.1	Méthodes subword	13
2.9.2	Domain mismatch	13
2.10	Limites et transition	14
3	PRETRAINING DATA: CORPORA AND KNOWLEDGE	15
3.1	Introduction	15
3.2	Web-Scale Pretraining Corpora	15
3.2.1	Common Crawl	15
3.2.2	Premiers grands corpus	16
3.2.3	Corpus modernes (2024-2025)	16
3.3	Data Quality and Curation	17
3.3.1	Filtrage heuristique	17
3.3.2	Classieurs de qualité	17
3.3.3	Organisation par domaine	18
3.3.4	Deduplication	18
3.3.5	Contamination	19
3.4	Biomedical and Clinical Corpora	19
3.4.1	Grandes sources anglophones	19
3.4.2	Curation article-level pour le biomed	20
3.4.3	Sources francophones	20
3.5	Medical Ontologies and Knowledge Graphs	21
3.5.1	Principales ontologies médicales	21
3.5.2	Représentations formelles	22
3.6	Knowledge-Enhanced Pretraining	23
3.6.1	Embeddings statiques de graphes	23
3.6.2	Injection de connaissances dans les transformers	23
3.6.3	Approches contrastives médicales	24
3.6.4	Graph-to-text pour le prétraining	24
3.7	Limites et transition	25
4	CLINICAL INFORMATION EXTRACTION	27
II	BUILDING A BIOMEDICAL CORPUS	29
5	COLLECTING BIOMEDICAL TEXT	31
6	FILTERING BY QUALITY SIGNALS	33
7	DETECTING CONTENT TYPES	35

III PRETRAINING LANGUAGE MODELS	37
8 ENCODER MODELS FOR FRENCH BIOMEDICINE	39
9 WHEN DECODER CONTINUE-PRETRAINING STOPS WORKING	41
10 BEYOND MASKED LANGUAGE MODELING	43
IV ADAPTING TO CLINICAL TASKS	45
11 LIMITS OF DIRECT FINE-TUNING	47
12 ARCHITECTURES FOR LOW-RESOURCE EXTRACTION	49
13 SYNTHETIC DATA FOR TASK ADAPTATION	51
14 CONCLUSION	53
ACRONYMS	55
GLOSSARY	57

1 INTRODUCTION

“The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin. [...] We should build in only the meta-methods that can find and capture this arbitrary complexity. [...] The eventual success is tinged with bitterness, and often incompletely digested.”

- Richard S. Sutton, *The Bitter Lesson*

Healthcare has been one of the earliest and most prominent application domains for artificial intelligence. Use cases range from diagnostic support to medical record management, cohort selection, report generation, and clinical information extraction.

As early as the 1970s, MYCIN, the iconic expert system developed at Stanford, diagnosed bacterial infections and recommended antibiotics. It already achieved performance comparable to human specialists in its narrow domain. INTERNIST-1, later known as CADUCEUS, tackled internal medicine diagnosis across hundreds of diseases. However, these systems relied on explicit rules and static knowledge bases, which limited their ability to handle the complexity and variability of real-world medical data. They were rigid and required constant maintenance to keep up with medical advances. Building such systems demanded expensive medical experts to annotate clinical data and write rules by hand.

With the rise of machine learning in the 1990s and 2000s, data-driven approaches began to dominate healthcare AI. Supervised and unsupervised learning algorithms trained models directly from data. The interpretability of rule-based systems was traded for the performance of deep neural networks capable of capturing complex patterns. R2 ImageChecker (1998) became the first FDA-approved computer-aided detection system for mammography, using neural networks to spot suspicious microcalcifications. PAPNET applied neural networks to cytological screening for cervical cancer.

Then came the scaling era. The recipe became simple: pretrain a language model on billions of words of raw text, then fine-tune on the target task. This approach routinely outperformed systems built with years of domain expertise. The need for expensive annotation, hand-crafted rules, and curated knowledge bases faded, beaten by data and compute.

Despite healthcare offering many structured and curated knowledge bases such as UMLS, SNOMED-CT, and RxNorm, experience has shown that general-purpose language models pre-trained on massive unstructured data outperform approaches based on domain-specific rules and knowledge bases. It seems our models cannot effectively integrate these structured resources, or perhaps these knowledge bases do not cover the full complexity and variability of medical knowledge as it is used in practice.

However, healthcare faces a limitation that other domains do not: the confidentiality and scarcity of annotated clinical data. Medical data is sensitive and difficult to obtain in large quantities for

1 Introduction

model training. Works adapting language models to the medical domain often rely on public data such as PubMed or patient forums, but these do not always reflect the real clinical language used in electronic health records. The clinical domain has unique jargon, abbreviations, and writing styles that differ from general medical texts. Moreover, while MIMIC exists in the United States, public clinical datasets remain scarce in other countries such as France.

This leads to a fundamental problem: the paradigm that has won is pretraining at scale, but in healthcare we lack massive public clinical data to pretrain language models. Can we find new ways to exploit public data for pretraining language models adapted to the clinical domain, without using sensitive clinical data?

PART I

RELATED WORKS

2 LANGUAGE MODELS

2.1 INTRODUCTION

Un language model c'est un modèle probabiliste sur des séquences de tokens:

$$P(\mathbf{w}) = \prod_{t=1}^L P(w_t | w_{<t})$$

Factorisation autoregressive (chain rule). Fondation de tout le NLP moderne.
Pourquoi c'est important pour la thèse:

- Clinical IE repose sur des LMs pretrained
- Domain adaptation = défi central
- French + biomedical + clinical = triple rareté de données

-> plan du chapitre: on remonte des n-grams jusqu'aux archis modernes, puis on zoomé sur le continual pretraining qui est au cœur de la thèse

2.2 FROM STATISTICAL TO NEURAL LANGUAGE MODELS

2.2.1 N-GRAM MODELS

Hypothèse de Markov: $P(w_t | w_{<t}) \approx P(w_t | w_{t-n+1}, \dots, w_{t-1})$.

- Kneser-Ney smoothing
- Problème: curse of dimensionality, $|\mathcal{V}|^n$ explose

2.2.2 NEURAL LMs

- Bengio et al. (2003): neural probabilistic LM
 - Embeddings denses, résout la sparsité
 - Mots similaires -> vecteurs proches
- Word2Vec (Mikolov et al., 2013): skip-gram / CBOW
 - $\vec{\text{king}} - \vec{\text{man}} + \vec{\text{woman}} \approx \vec{\text{queen}}$
- ELMo (Peters et al., 2018): embeddings contextuels via biLSTM

2 Language Models

- Meme mot, contextes différents -> vecteurs différents
- Precursor du pretrain-then-finetune

2.2.3 RNNs

- LSTM (Hochreiter & Schmidhuber, 1997)
 - En théorie: long-range dependencies
 - En pratique: vanishing gradients, pas parallelisable
- > bottleneck pour scaler -> besoin d'une nouvelle archi

2.3 TRANSFORMER

2.3.1 SELF-ATTENTION

Vaswani et al. (2017).

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right)V$$

- Interactions directes cross-position (vs RNN indirect via hidden state)
- $O(L^2)$ temps et mémoire
- Parallelisable
- Multi-head: h têtes, chacune spécialise (syntaxe, sémantique, position)

2.3.2 POSITIONAL ENCODINGS

- APE sinusoidal (Vaswani 2017) -> extrapolation limitée
- APE learned (Devlin 2019) -> pas d'extrapolation
- ALiBi (Press et al. 2022) -> bias linéaire, mieux
- RoPE (Su et al. 2021) -> matrices de rotation, le meilleur

RoPE = standard maintenant: LLaMA, Mistral, ModernBERT.

2.3.3 ENCODER VS DECODER

- Encoder: attention bidirectionnelle, MLM, classification/NER (BERT, RoBERTa)
 - Decoder: attention causale, CLM, génération (GPT, LLaMA)
 - Encoder-decoder: T5 (Raffel et al., 2020), seq2seq, plus trop utilisé
- > maintenant qu'on a l'archi, question: quel objectif de prétraining?

2.4 PRETRAINING OBJECTIVES

2.4.1 CLM (CAUSAL LANGUAGE MODELING)

GPT (Radford et al., 2018).

$$\mathcal{L}_{\text{CLM}} = - \sum_{t=1}^L \log P(w_t | w_{<t})$$

- Tous les tokens contribuent au gradient
- Supervision dense
- Naturel pour la generation

2.4.2 MLM (MASKED LANGUAGE MODELING)

BERT (Devlin et al., 2019).

$$\mathcal{L}_{\text{MLM}} = - \sum_{t \in \mathcal{M}} \log P(w_t | \mathbf{w}_{\setminus t})$$

- 15% tokens masked (80% [MASK], 10% random, 10% unchanged)
- Contexte bidirectionnel -> representations plus riches
- Mais: seulement 15% contribuent -> gradient sparse
- Contraste avec CLM ou chaque token compte

2.4.3 OBJECTIFS HYBRIDES (2024-2025)

Important pour le Ch.5 de la these.

Approche CLM puis MLM:

- CLM d'abord puis decay vers MLM = meilleur que MLM seul
- CLM: gradients denses (tous les tokens), MLM: sparse (15%)
- Hypothese: CLM “compresse” la connaissance, MLM restaure la bidirectionnalite

Refs:

- “Should We Still Pretrain Encoders with MLM?” (arXiv:2507.00994, 2025) -> biphasic CLM vers MLM bat pure MLM
- GPT-BERT (arXiv:2410.24159, 2024) -> unifie CLM+MLM
- AntLM (CoNLL 2024) -> alternance CLM/MLM

Note: ELECTRA (Clark et al., 2020) avait propose RTD pour signal plus dense, mais Modern-BERT montre que bonne archi + MLM suffit (bat DeBERTaV3 qui utilise RTD).

-> on a les archis et les objectifs, question maintenant: est-ce que ca scale?

2 Language Models

2.5 SCALING LAWS & LLMs

2.5.1 SCALING

- GPT-3 (Brown et al., 2020): 175B params, few-shot emergent, prompting remplace le finetuning
- Chinchilla (Hoffmann et al., 2022):

$$L(N, D) = \frac{A}{N^\alpha} + \frac{B}{D^\beta} + E$$

-> compute-optimal: smaller model + more data bat larger model + less data

- LLaMA (Touvron et al., 2023): 7B rivalise avec bcp plus gros via 1T tokens -> open weights, democratisation

2.5.2 POST-CHATGPT

- RLHF (Ouyang et al., 2022) -> preference learning
- Instruction tuning
- ChatGPT, Claude, Gemini
- Emergent abilities (Wei et al., 2022) -> chain-of-thought etc.

-> ok ca scale pour le general, mais nous on a un domaine specifique -> comment adapter?

2.6 CONTINUAL PRETRAINING

Section centrale pour la these. C'est le coeur de ce qu'on fait (CamemBERT-bio, ModernCamemBERT-bio).

2.6.1 FONDEMENTS ET TERMINOLOGIE

“Don’t Stop Pretraining” (Gururangan et al., ACL 2020):

- DAPT = Domain-Adaptive Pre-Training: 2eme phase sur corpus domaine
- TAPT = Task-Adaptive Pre-Training: sur donnees non-annotatees de la tache
- DAPT + TAPT combine -> meilleurs resultats
- Jusqu'a +3 F1 meme avec petit corpus

Taxonomie moderne (Wang et al., CSUR 2025):

- CPT = nouvelles donnees generales au fil du temps
- DAP = specialisation domaine
- CFT = adaptation sequentielle de taches

2.6.2 CATASTROPHIC FORGETTING

Le probleme (Kirkpatrick et al., 2017):

- Les reseaux ecrasent les poids necessaires aux taches precedentes
- EWC (Elastic Weight Consolidation): penaliser les changements sur les poids importants

Pourquoi ca compte pour les LLMs:

- Continual PT peut degrader les capacites generales
- Capacite a suivre des instructions particulierement fragile
- Petits modeles plus sensibles (Yildiz et al., 2024)

“Spurious forgetting” (2024): les baisses de perf refletent peut-être une perte d’alignement, pas une vraie perte de connaissance.

-> ok c'est un vrai probleme, quelles solutions?

2.6.3 SOLUTIONS RECENTES (2024-2025)

Ibrahim et al. (2024), “Simple and Scalable Strategies”:

- LR re-warming: relancer le learning rate quand on ajoute des donnees
- LR re-decaying: cosine decay apres warmup
- Replay: mixer 1% des donnees precedentes -> ca suffit
- Valide a 405M et 10B
- -> egalise from-scratch pour une fraction du compute

Stability gap (arXiv:2406.14833, 2024):

- Chute temporaire au début du continual PT, puis recovery
- Solutions: subset propre, sous-corpus de qualite, data mixing
- 36.2% -> 40.7% avec seulement 40% du budget training

Data selection:

- Donnees non pertinentes = degradation
- 10% du corpus bien selectionne fait aussi bien que 100% vanilla continual PT

2 Language Models

2.6.4 FROM SCRATCH VS CONTINUAL PT

- From scratch (PubMedBERT, DrBERT): vocab custom, pas de forgetting, mais tres cher
- Continual (BioBERT, Gururangan): efficient, preserve le general, mais tokenizer sous-optimal

Cout compute:

- DrBERT from scratch: 128x V100, 20h
- AliBERT from scratch: 48x A100, 20h
- BioBERT continual: 8x V100, ~10 jours
- -> from-scratch coutre ~30x plus de compute

2.6.5 ENCODERS VS DECODERS: DYNAMIQUES DIFFERENTES

Decoders (LLMs), ca marche pas bien:

- BioMistral: -0.9 points apres continual PT
- Dorfner et al. (2024): biomedical LLMs sous-performent les generalistes
- OpenBioLLM-8B: 30% vs Llama-3-8B: 64.3% sur cas NEJM
- -> hypothese: perte de connaissances generales pendant la specialisation

Encoders, c'est different:

- Contexte bidirectionnel peut-etre plus robuste au forgetting
- MLM moins sujet au catastrophic forgetting?
- -> voir Ch.5: approche CLM vers MLM sur les encoders

-> tension narrative: continual PT marche pas pour les decoders, et pour les encoders? -> c'est la question de la these

2.6.6 TENTATIVES ANTERIEURES EN FRANCAIS BIOMEDICAL

- Copara et al. (2020): 31K articles -> aucune amelioration (corpus trop petit)
- Le Clercq de Lannoy et al. (2022): 136M mots -> +2 EMEA seulement
- Dura et al. (2022): 21M docs APHP -> +3% mais donnees privees

Question ouverte: pourquoi des resultats contradictoires? DrBERT dit que continual PT marche pas. CamemBERT-bio (cette these) montre que si, avec le bon corpus.

-> ok on sait ce qu'est le continual PT, maintenant quels modeles existent?

2.7 MODELES PRETRAINED BIOMED

2.7.1 ENCODERS BIOMED ANGLAIS

- BioBERT (Lee et al., 2019): PubMed+PMC 18B mots, continual -> +0.62% NER, +2.80% RE
- SciBERT (Beltagy et al., 2019): Semantic Scholar 1.14M papers, from scratch -> 42% vocab overlap avec BERT
- PubMedBERT (Gu et al., 2022): PubMed only, from scratch -> introduit BLURB benchmark
- ClinicalBERT (Alsentzer et al., 2019): MIMIC-III, continual -> acces restreint
- BioLinkBERT (Yasunaga et al., 2022): PubMed + citation links, from scratch -> +7% BioASQ

SOTA 2025: BioClinical ModernBERT, cf section architectures modernes.

2.7.2 MODELES FRANCAIS

General:

- CamemBERT (Martin et al., 2020): RoBERTa sur OSCAR 138GB -> resultat interessant: 4GB donne a peu pres la meme perf que 138GB
- FlauBERT (Le et al., 2020)
- CamemBERT 2.0 / CamemBERTa (2024): architecture DeBERTaV3, nouveau tokenizer

Biomedical francais:

- DrBERT (Labrak et al., 2023): from scratch, public + prive -> pretend que continual PT marche pas pour le francais biomed
- AlIBERT (Berhe et al., 2023): from scratch, Unigram regularise -> pas disponible publiquement

-> tension: DrBERT dit from-scratch necessaire, CamemBERT-bio (cette these) montre le contraire

2.7.3 DECODERS BIOMED

- BioMistral (Labrak et al., 2024): 3B tokens PMC, 32x A100 20h, base Mistral
- Meditron (Chen et al., 2023): 46B tokens, 128x A100 332h, base Llama-2
- PMC-LLaMA (Wu et al., 2024): 75B tokens, massif

Continual PT pour decoders = resultats mitiges:

2 Language Models

- BioMistral: -0.9 points, recupere seulement via model merging
 - Dorfner et al. (2024): “Biomedical LLMs Seem not to be Superior to Generalist Models”
 - -> contraste avec les encoders, voir Ch.5
- > les modeles existent mais les archis ont evolue, quoi de neuf?

2.8 ARCHITECTURES MODERNES

2.8.1 MODERNISATION DES ENCODERS

DeBERTa (He et al., ICLR 2021):

- Disentangled attention: separe contenu et position
- Premier a battre l'humain sur SuperGLUE
- -> influence sur ModernBERT, CamemBERT 2.0

MosaicBERT (Portes et al., NeurIPS 2023):

- 30% masking au lieu de 15% BERT
- FlashAttention + ALiBi + GLU + unpadding
- BERT-base en 1.13h sur 8x A100 (~\$20)
- -> fondation directe de ModernBERT

2.8.2 EFFICIENT ATTENTION

FlashAttention (Dao et al., 2022):

- Reorganise le calcul pour la hierarchie memoire GPU (tiling, kernel fusion)
- Meme $O(L^2)$ mais 2-4x plus rapide en pratique
- Standard maintenant

Autres:

- GQA (Ainslie et al., 2023): partage de KV-heads -> Llama-2, Mistral
- RMSNorm (Zhang & Sennrich, 2019): LayerNorm simplifie

2.8.3 LONG CONTEXT

Le problème: BERT = 512 tokens, comptes-rendus cliniques = souvent 2000+. Tronquer = perdre de l'info diagnostique.

- Longformer (Beltagy et al., 2020): 4096 tokens, sliding window + global attention
- BigBird (Zaheer et al., 2020): 4096, sparse attention
- Clinical Longformer (Li et al., 2022): 4096, Longformer continue pretrained sur MIMIC-III
- ModernBERT (Warner et al., 2024): 8192 tokens
 - FlashAttention v2/v3
 - RoPE
 - Alternance local/global attention
 - 30% masking (de MosaicBERT)
 - -> 16x le contexte de BERT

Adaptations cliniques 2025:

- Clinical ModernBERT (Lee et al., 2025): PubMed + MIMIC-IV + ontologies
- BioClinical ModernBERT (2025): 53.5B tokens, 20 datasets cliniques -> SOTA actuel

2.9 TOKENIZATION

2.9.1 METHODES SUBWORD

- BPE (Sennrich et al., 2016): merge les paires les plus fréquentes
- SentencePiece (Kudo & Richardson, 2018): language-independent
- Unigram (Kudo, 2018): sélection probabiliste

2.9.2 DOMAIN MISMATCH

SciBERT: vocab général vs scientifique = 42% intersection seulement.

Exemple biomed français:

- Tokenizer général: "echocardiographie" -> ["echo", "#cardi", "#graphie"]
- Tokenizer domaine: "echocardiographie" -> ["echocardiographi", "e"]

Tradeoff:

- From-scratch = vocab custom mais cher
- Continual PT = tokenizer sous-optimal mais efficient

2.10 LIMITES ET TRANSITION

- Rarete des donnees: texte clinique prive (CNIL), modeles hospitaliers pas partageables -> besoin d'alternatives publiques -> Ch.3 Corpus Annotation
- Knowledge gap: LMs captent des patterns, pas la structure des ontologies. Codage medical = connaissances hierarchiques -> approches knowledge-enhanced -> Ch.4 Clinical IE
- Documents longs: comptes-rendus hospitaliers depassent les limites de BERT. Longformer et ModernBERT aident mais faut adapter au domaine

-> ce chapitre a couvert les fondations techniques du language modeling. Le chapitre suivant traite de comment construire et annoter des corpus pour entrainer ces modeles.

3

PRETRAINING DATA: CORPORA AND KNOWLEDGE

3.1 INTRODUCTION

Le chapitre précédent couvrait les modèles. Maintenant faut les nourrir.

Les LMs modernes sont data-hungry:

- GPT-3 (Brown et al., 2020): 300B tokens
- LLaMA (Touvron et al., 2023): 1-2T tokens
- Chinchilla (Hoffmann et al., 2022): compute-optimal = plus de data, modèle plus petit

La qualité et la composition du corpus de prétraining déterminent largement les capacités du modèle. "Data is the new bottleneck."

Deux types de données de prétraining:

- Texte brut: web crawls, articles scientifiques, livres
- Connaissances structurees: ontologies médicales, knowledge graphs

Pourquoi c'est central pour la thèse:

- Part 1 = construire un corpus biomédical français à partir de sources publiques
- Ch.6 = injecter des connaissances d'ontologies médicales pendant le prétraining
- Triple rareté: français + biomedical + clinique

-> plan du chapitre: on part des corpus web généralistes, on zoom sur le filtrage qualité, puis sur les sources biomédicales spécifiques, et enfin sur les ontologies médicales et comment les intégrer dans le prétraining

3.2 WEB-SCALE PRETRAINING CORPORA

3.2.1 COMMON CRAWL

Source de base de presque tous les grands corpus de prétraining:

- Web crawl ouvert, 2.5B pages par snapshot mensuel, 250-400 TiB non compressé
- Snapshots mensuels depuis 2008, petaoctets cumulés
- Contenu très hétérogène: articles, forums, spam, code, pubs, boilerplate...
- Pas utilisable directement -> nécessite des pipelines de nettoyage lourdes

3 Pretraining Data: Corpora and Knowledge

3.2.2 PREMIERS GRANDS CORPUS

C4 (Raffel et al., 2020):

- Colossal Clean Crawled Corpus, 750GB d'anglais extrait d'un seul snapshot CC (avril 2019)
- Filtres: langdetect $\geq 99\%$ anglais, suppression phrases incomplètes, dedup
- Dodge et al. (2021) documentent le contenu: brevets, sites militaires US, texte machine-generated, exemples de benchmarks
- Blocklist filtering retire disproportionnellement du texte de/sur les minorités

The Pile (Gao et al., 2020):

- 800GB anglais, 22 sources diversifiées (PubMed, ArXiv, GitHub, StackExchange, Wikipedia, livres, etc.)
- Idée: diversité des sources > volume brut d'une seule source
- Open-source (EleutherAI)

OSCAR (Ortiz Suarez et al., 2019):

- Pipeline asynchrone pour classifier CC par langue
- Optimisé pour infra moyenne/basse ressource (I/O-bound)
- Produit le corpus de CamemBERT: 138GB de texte français
- Shuffled au niveau ligne pour éviter les problèmes de copyright (ironiquement)

3.2.3 CORPUS MODERNES (2024-2025)

FineWeb (Penedo et al., 2024):

- 15T tokens, 96 snapshots CC
- Documentation détaillée des stratégies de dedup et filtrage
- Pipeline: extraction texte, filtrage langue, heuristiques qualité, dedup URL + MinHash
- FineWeb-Edu: sous-ensemble 1.3T tokens filtre pour contenu éducatif (cf section suivante)

RedPajama-V2 (Weber et al., 2024):

- Dataset ouvert pour entraîner des LLMs
- Inclut subset français (RPv2-Fr)
- Métriques de qualité pré-calculées pour faciliter le filtrage

TxT360 (Tang et al., 2024):

- "Perfect blend" de sources
- Dedup globale sur 99 snapshots CC: 20TB -> 4.83T tokens (reduction 80%)
- Pipeline similaire a FineWeb + near-dedup globale additionnelle

Autres:

- Dolma (Soldaini et al., 2024): 3T tokens, corpus ouvert pour recherche (OLMo)
- ROOTS (Laurencon et al., 2023): 1.6TB, multilingue composite (BigScience/BLOOM)

-> on a des milliards de tokens de texte web, mais la qualite varie enormement -> comment selectionner les bons documents?

3.3 DATA QUALITY AND CURATION

3.3.1 FILTRAGE HEURISTIQUE

Approches classiques, devenues standard depuis Gopher (Rae et al., 2022):

- Regles: longueur min/max du document, ratio mots/symboles, proportion de majuscules, lignes dupliquees
- Filtrage par perplexite: KenLM entraîne sur Wikipedia, documents à perplexité trop haute = bruit
- Language ID: fastText pour filtrer par langue
- Ces heuristiques sont reprises par presque tous les pipelines modernes (FineWeb, C4, Dolma)

Limites: les heuristiques sont nécessaires mais insuffisantes. Elles retirent le bruit évident mais ne distinguent pas la qualité du contenu.

3.3.2 CLASSIFIEURS DE QUALITE

Paradigme récent: LLM-as-annotator -> distillation vers petit classifieur -> filtrage à grande échelle. FineWeb-Edu (Penedo et al., 2024):

- Llama-3-70B-Instruct annote 460K pages web sur "educational value" (échelle 0-5)
- Échelle additive: le LLM évalue chaque critère et construit le score pas à pas
- Prompt cible le niveau primaire/collège pour éviter de favoriser les papiers trop techniques
- Distillation vers classifieur de régression (Snowflake-arctic-embed, F1 = 82% en binaire avec seuil >= 3)
- Classification de 15T tokens: 6K heures H100

3 Pretraining Data: Corpora and Knowledge

- Retire 92% des données, bat quand même le corpus complet sur MMLU, ARC, Open-BookQA

DCLM / DataComp-LM (Li et al., 2024):

- Classifieur FastText entraîné pour séparer instructions synthétiques (Open Hermes 2.5) du web général
- Favorise les structures Q&A résolues: question courte + réponse directe + raisonnement bref
- Utilisé dans le mix de prétraining d'OLMo2

Nemotron-CC (NVIDIA, 2024):

- Classifieur de qualité multi-critères (accuracy, clarity, coherence, etc.)
- Approche plus granulaire que le score éducatif unique de FineWeb-Edu

3.3.3 ORGANISATION PAR DOMAINE

WebOrganizer (Wettig et al., 2025):

- Organise le web en taxonomies par topic et par format
- Annotations LLM distillées vers classificateurs efficaces
- Montre que FineWeb-Edu a un biais implicite vers Science & Technology et Academic Writing
- L'efficacité de FineWeb-Edu vient en partie de préférences de domaine alignées avec les benchmarks (MMLU, HellaSwag)
- Combiner organisation par domaine + filtrage qualité > qualité seule

3.3.4 DEDUPLICATION

- URL-level: retirer les duplicates exacts d'URL
- Document-level: MinHash + LSH pour near-duplicates
- TxT360: dedup globale sur 99 snapshots CC, réduction de 80%
- Crucial: la duplication biaise le training vers les textes les plus copies (boilerplate, templates)

3.3.5 CONTAMINATION

Le filtrage agressif pose un probleme de contamination des benchmarks:

- Deng et al. (2024): ChatGPT et GPT-4 devinent 52% et 57% des options manquantes de MMLU -> les QA pairs sont sur le web
- Xu et al. (2024): survey des techniques de detection de contamination dans les LLMs
- InfiniGram (Liu et al., 2024): outil pour identifier des exact matches dans les training sets

Tension fondamentale: filtrer plus agressivement = meilleurs scores benchmarks, mais potentiellement via contamination implicite plutot que vraie qualite.

-> ok le filtrage marche pour les corpus web generalistes, mais en biomedical c'est different: les corpus sont plus petits, le domaine est tres specifique, les criteres de qualite ne sont pas les memes
-> quelles sources biomed existent?

3.4 BIOMEDICAL AND CLINICAL CORPORA

3.4.1 GRANDES SOURCES ANGLOPHONES

PubMed:

- Base de donnees de la NLM (National Library of Medicine)
- 39M citations et abstracts d'articles biomedicaux
- Texte structure, vocabulaire technique, qualite editoriale
- Abstracts seulement (pas le full-text)
- Utilise par BioBERT (Lee et al., 2019), PubMedBERT (Gu et al., 2022)

PMC Open Access:

- 4.5M articles full-text en acces libre
- 98% anglais (Labrak et al., 2024)
- Contenu heterogene: articles de recherche, reviews, case reports, editoriaux, guidelines
- Utilise par BioMistral (3B tokens), Meditron (46B tokens), PMC-LLaMA (75B tokens)
- Probleme: granularite article. Un article de recherche peut contenir a la fois des cas cliniques pertinents et des sections methodologiques non pertinentes

MIMIC (Johnson et al., 2016/2023):

- MIMIC-III, MIMIC-IV: notes cliniques du Beth Israel Deaconess Medical Center
- Acces restreint (data use agreement, formation ethique)

3 Pretraining Data: Corpora and Knowledge

- Utilise par ClinicalBERT (Alsentzer et al., 2019), Clinical Longformer (Li et al., 2022), BioClinical ModernBERT (Sounack et al., 2025)
- Gold standard pour le texte clinique anglais, mais un seul hopital americain -> pas generalizable

3.4.2 CURATION ARTICLE-LEVEL POUR LE BIOMED

Meditron (Chen et al., 2023):

- Score 0-1 par article base sur: MeSH tags, type de publication, reputation du journal, recence, nombre de citations
- Upsampling des articles a haut score, filtrage des bas scores
- 46B tokens au total, 128x A100, 332h pour le 70B
- Limite: granularite article. Un article peut melanger contenu pertinent et non-pertinent

BioClinical ModernBERT (Sounack et al., 2025):

- 169B tokens: PubMed + PMC + 20 datasets cliniques (MIMIC-III, MIMIC-IV, CheXpert Plus, etc.)
- SOTA actuel sur les benchmarks biomed anglais
- Mais: la plupart des sources cliniques necessitent des data use agreements -> pas vraiment "public"

3.4.3 SOURCES FRANCOPHONES

Ressources publiques existantes:

- ISTELEX: base de 27M references scientifiques, permet d'extraire des documents francais bio/med
- CLEAR (Grabar & Cardon, 2018): notices medicamenteuses, articles d'encyclopedie medcale en francais technique et simplifie
- E3C (Magnini et al., 2021): corpus multilingue de cas cliniques, notices, resumes de these medicale
- HAL / Halvest (Kulumba et al., 2024): articles scientifiques francais et anglais extraits du depot HAL (Hyper Articles en Ligne)
- Wikipedia biomedical: articles des portails medecine, pharmacie, biologie. Libre et multi-lingue

Modeles entraines from scratch sur corpus francais:

- DrBERT (Labrak et al., 2023): corpus de 7GB, from scratch sur 128x V100

- Pretend que le continual PT ne marche pas pour le francais biomed
- Mais: evaluation potentiellement biaisee (cf Ch.4 de la these)
- AliBERT (Berhe et al., 2023): ScienceDirect + theses Sudoc, tokenizer Unigram regularise, 48x A100
 - Modele non disponible publiquement

Tentatives de continual pretraining en francais biomed:

- Copara et al. (2020): 31K articles seulement -> aucune amelioration sur CamemBERT-large
- Le Clercq de Lannoy et al. (2022): PubMed + Cochrane + ISTEY + Wikipedia = 136M mots -> +2 F1 sur EMEA, rien sur MEDLINE
- Dura et al. (2022): 21M documents cliniques APHP -> +3% sur APMed (privé), mais donnees non partageables

Resultats décevants: DrBERT dit que continual PT marche pas, Copara et Le Clercq confirment. Mais le probleme est peut-etre le corpus (trop petit, pas assez cible), pas la methode.

Triple rarete: francais + biomedical + clinique. Les textes cliniques sont privés (CNIL), les modèles hospitaliers ne sont pas partageables entre établissements. Les sources publiques francophones sont fragmentées et petites comparées à PubMed/PMC.

-> le texte biomedical aide pour le pretraining, mais les LMs captent des patterns statistiques, pas la structure des ontologies. Pour le codage medical (CIM-10, CCAM) il faut des connaissances hiérarchiques explicites -> peut-on utiliser directement les ontologies comme données de pretraining?

3.5 MEDICAL ONTOLOGIES AND KNOWLEDGE GRAPHS

3.5.1 PRINCIPALES ONTOLOGIES MÉDICALES

UMLS (Lindberg et al., 1993):

- Unified Medical Language System, meta-ontologie biomedical de la NLM
- Intègre 189 vocabulaires sources (SNOMED-CT, MeSH, ICD, ATC, etc.)
- 3.4M concepts, 16.7M noms de concepts uniques
- Standard international pour la recherche en NLP biomédical
- Contient: synonymes, relations hiérarchiques, relations associatives, définitions

SNOMED-CT:

- 371K concepts cliniques structurés
- Relations hiérarchiques (is-a) + associatives (finding site, causative agent, etc.)

3 Pretraining Data: Corpora and Knowledge

- Standard pour les systèmes d'information clinique internationaux
- Couverture: diagnostics, procédures, anatomie, substances, organismes

CIM-10 / ICD-10 (OMS, 2019):

- Classification internationale des maladies
- Hierarchy: chapitres -> blocs -> categories -> sous-categories
- Version française enrichie par l'ATIH avec: exclusions, inclusions, notes, libelles courts/longs
- Utilisé pour le codage PMSI (Programme de Médicalisation des Systèmes d'Information) en France
- 19K codes dans la version ATIH

CCAM (ATIH):

- Classification Commune des Actes Médicaux, 8K codes
- Spécifique à la France, structure hiérarchique
- Relations: actes associés, incompatibilités, modificateurs

ATC (OMS):

- Anatomical Therapeutic Chemical, 6K codes au 5ème niveau
- 5 niveaux: anatomique -> thérapeutique -> pharmacologique -> chimique -> substance
- Standard pour la classification des médicaments

3.5.2 REPRESENTATIONS FORMELLES

- RDF (Resource Description Framework): triplets sujet-prédicat-objet, standard W3C
- OWL (Web Ontology Language): logique descriptive, raisonnement automatique, classes, propriétés
- SKOS (Simple Knowledge Organization System): prefLabel, altLabel (synonymes), broader/narrower (hierarchy), related, exactMatch

Les ontologies codent explicitement ce que les LMs doivent apprendre implicitement:

- Hierarchy: E11 (diabète type 2) ⊂ E10-E14 (diabète sucre)
- Exclusions: E11 ≠ E10 (diabète type 1) – ne pas confondre
- Synonyms: un même concept a plusieurs libellés (prefLabel + altLabel)
- Relations associatives: complications, traitements, liens entre ontologies

-> les ontologies sont riches en connaissances structurées, mais sous forme de graphes: pas directement exploitables par les LMs qui attendent du texte en langage naturel -> comment les intégrer dans le prétraining?

3.6 KNOWLEDGE-ENHANCED PRETRAINING

3.6.1 EMBEDDINGS STATIQUES DE GRAPHES

Approches qui produisent des vecteurs fixes (pas contextuels) a partir de graphes:

- Node2Vec (Grover & Leskovec, 2016): random walks sur graphe -> skip-gram, apprentissage de representations de noeuds
- RDF2Vec (Ristoski & Paulheim, 2016): adapte le paradigme walk + skip-gram aux graphes RDF. Fondation des approches de walks sur ontologies
- Snomed2Vec (Agarwal et al., 2019): Node2Vec applique a SNOMED-CT, 5-6x amelioration sur concept similarity
- OWL2Vec* (Chen et al., 2021): walks sur ontologies OWL (teste sur Gene Ontology)

Limite: embeddings statiques, pas integres dans un transformer. Un concept = un vecteur fixe quel que soit le contexte. Pas de contextualisation.

-> on peut faire mieux en integrant les KG directement dans les transformers

3.6.2 INJECTION DE CONNAISSANCES DANS LES TRANSFORMERS

Premiere generation: injecter des entites dans l'architecture:

- ERNIE (Zhang et al., 2019): aligne entity embeddings (TransE) avec representations textuelles pendant le pretraining. Entity linking implicite
- KnowBERT (Peters et al., 2019): integre des entity linkers comme couches d'attention supplementaires dans BERT
- K-BERT (Liu et al., 2020): injecte des triples de KG comme soft-position embeddings dans l'input. Modifie la sequence d'entree directement

Deuxieme generation: multi-task pretraining language + KG:

- KEPLER (Wang et al., 2021): MLM + TransE knowledge embedding sur Wikidata simultanement. Les deux objectifs partagent le meme encoder
- CoLAKE (Sun et al., COLING 2020): construit des word-knowledge graphs qui unifient contexte textuel et contexte KG, puis pretraigne avec MLM modifie sur ces graphes
- DRAGON (Yasunaga et al., NeurIPS 2022): joint pretraining bidirectionnel MLM + link prediction sur UMLS
 - PubMed pour le texte, UMLS pour le KG
 - +3% accuracy sur MedQA, +10% sur raisonnement multi-step
 - Reference pour le multi-task MLM + KG

3 Pretraining Data: Corpora and Knowledge

3.6.3 APPROCHES CONTRASTIVES MEDICALES

- SapBERT (Liu et al., NAACL 2021): self-alignment pretraining sur synonymes UMLS
 - Contrastive learning: synonymes UMLS = positifs, non-synonymes = negatifs
 - Metric learning scalable sur 4M+ concepts UMLS
 - SOTA entity linking biomed sur 6 benchmarks
 - Un seul modèle pour tous les types d'entités ("one-model-for-all")
- CODER (Yuan et al., 2022): contrastive pretraining sur 87M triplets de relations UMLS
 - Relations + synonymes dans l'objectif contrastif
 - Medical term normalization

3.6.4 GRAPH-TO-TEXT POUR LE PRETRAINING

Approche récente: transformer les graphes en texte lisible par les LMs.

"Textbooks Are All You Need" (Gunasekar et al., 2023):

- Données synthétiques de qualité textbook générées par GPT -> petits modèles très performants
- Montre que la qualité et la structure des données importent plus que le volume brut
- Paradigme: générer des données synthétiques structurées plutôt que crawler du texte brut

AntGLM-Med (Zhou et al., 2023):

- Recrit des sous-graphes de KG medical en prose naturelle
- Utilise pour le continual pretraining de LLMs médicaux
- Plus proche précédent de la génération de texte à partir d'ontologies pour pretraining

EntiGraph (Yang et al., ICLR 2025):

- Algorithme d'augmentation entity-centric: décompose un corpus en entités puis génère du texte sur les relations entre entités
- 1.3M tokens réels -> 600M tokens synthétiques via GPT-4-turbo
- Scaling log-linéaire de l'accuracy avec le nombre de tokens synthétiques
- Cible les décodeurs (Llama 3 8B), pas encore testé sur encoders

Limites actuelles du knowledge-enhanced pretraining:

- La plupart des travaux sont en anglais et sur UMLS
- Graph-to-text pour le pretraining est récent (2023-2025)

- Pas encore explore pour les encoders bidirectionnels (DRAGON fait du link prediction, pas du graph-to-text)
- Les ontologies nationales (CIM-10 ATIH, CCAM) n'ont jamais ete utilisees pour du pre-training
- SapBERT et CODER exploitent UMLS mais uniquement via contrastive learning, pas via generation de texte

-> methodes prometteuses mais gap important: langues non-anglaises, ontologies nationales, pretraining d'encoders via graph-to-text

3.7 LIMITES ET TRANSITION

- Rarete biomed non-anglais: PMC = 98% anglais, les corpus francais biomed sont petits et fragmentes. Les pipelines de filtrage (FineWeb-Edu, DCLM) sont concus pour l'anglais -> besoin de strategies de collection et de filtrage adaptees au domaine et a la langue
- Granularite de filtrage: les approches existantes (Meditron, FineWeb-Edu) filtrent au niveau document ou article. Mais un article biomed peut contenir a la fois des cas cliniques pertinents et des sections non pertinentes -> besoin de filtrage a granularite plus fine (paragraphe)
- Criteres de qualite non adaptes: les criteres "educatifs" de FineWeb-Edu ne correspondent pas au texte biomedical (cas cliniques, notices medicamenteuses, protocoles). Les biais vers Science & Technology identifies par WebOrganizer peuvent etre differents en biomed -> besoin de criteres domaine-specifiques
- Ontologies nationales sous-exploitees: UMLS domine la recherche, mais les systemes hospitaliers francais utilisent CIM-10/ATIH, CCAM, ATC. Aucun travail n'a integre ces ontologies dans le pretraining -> gap pour le codage medical en francais
- Graph-to-text naissant: AntGLM-Med et EntiGraph montrent la voie mais sont limites aux decoders en anglais. Pas encore explore pour les encoders bidirectionnels ni pour les ontologies non-anglophones -> possibilite de generer du texte "textbook" a partir des ontologies ATIH

-> ce chapitre a couvert les donnees de pretraining: sources textuelles web et biomed, methodes de curation et filtrage, ontologies medicales et leurs integrations dans le pretraining. Le chapitre suivant traite de la tache finale: l'extraction d'information clinique, et de comment les modeles pretrained sont utilises pour cette tache.

4

CLINICAL INFORMATION EXTRACTION

PART II

BUILDING A BIOMEDICAL CORPUS

5

COLLECTING BIOMEDICAL TEXT

6 FILTERING BY QUALITY SIGNALS

7

DETECTING CONTENT TYPES

PART III

PRETRAINING LANGUAGE MODELS

8

ENCODER MODELS FOR FRENCH BIOMEDICINE

9 WHEN DECODER CONTINUE-PRETRAINING STOPS WORKING

10 BEYOND MASKED LANGUAGE MODELING

PART IV

ADAPTING TO CLINICAL TASKS

11 LIMITS OF DIRECT FINE-TUNING

12 ARCHITECTURES FOR LOW-RESOURCE EXTRACTION

13 SYNTHETIC DATA FOR TASK ADAPTATION

14 CONCLUSION

ACRONYMS

PCA	Principal component analysis
SNF	Smith normal form
TDA	Topological data analysis

GLOSSARY

\LaTeX A document preparation system
 \mathbb{R} The set of real numbers

APPENDIX

