

Establishing clinical NLP modelling recommendations for restricted data availability settings

Anonymous ACL submission

Abstract

When solving clinical Natural Language Processing (NLP) downstream tasks, it is well-established that incorporating clinical-specific knowledge enhances model performances. However, there are scenarios where access to data or domain-specific models is not feasible. Despite various paradigms for adapting clinical NLP-based models, such as fine-tuning already pre-trained language models, pre-training and fine-tuning models, or in-context learning, the advantages of each alternative regarding data availability still need to be explored. We determined the impact of data availability and paradigm selection in the performance of models on solving multiple clinical NLP tasks in Spanish by simulating multiple clinical data availability settings and testing various NLP modelling paradigms. Overall, the best-performing modelling strategy was pre-training a masked language model (LM) with environment-specific unannotated text starting from an off-the-shelf clinical checkpoint and then fine-tuning the LM for the downstream task. The increase in performance from the continuation of pre-training of an off-the-shelf LM is marginal, considering the high amount of resources needed for the pre-training; therefore, we recommend the fine-tuning of an off-the-shelf clinical-specific LM if the model and task-specific data are available. We recommend a few-shot learning technique using a large LM if no data is available.

1 Introduction

Natural language processing (NLP) has gained tremendous importance in recent years with the advent of Transformer-based pre-trained language models (PLM) (Qiu et al., 2020) and large language models (LLM) (Zhao et al., 2023). These LMs have become the new paradigm for NLP-based machine learning modelling because of their modularity and ease of transferring learning. One can fine-tune

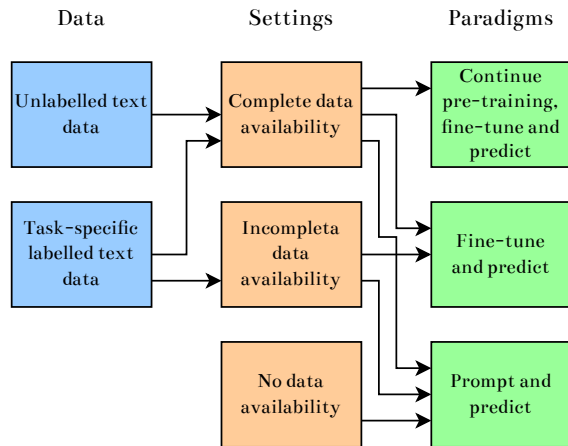
a PLM to solve any NLP task using off-the-shelf, already pre-trained LMs (Dodge et al., 2020).

It is known that using closer-to-the-domain LMs for fine-tuning downstream models improves the performance of the fine-tuned model (Gu et al., 2021; Zheng et al., 2022; Carrino et al., 2021). One of the most widespread paradigms for NLP-based modelling is the usage of a PLM for representing documents and then using a couple of layers to adapt the PLM output to solve the specific task; this framework is called the fine-tuning of PLMs. There are multiple options to optimize a model using this framework. The first option is to use a PLM and then fine-tune it to the downstream task or to pre-train an LM and then fine-tune it to the downstream task by employing unlabeled and task-labeled data. The second option is useful when no pre-trained models are available or one wants a closer-to-the-domain PLM. This second method involves initial pre-training of an LM utilizing unlabeled data specific to the target domain, followed by fine-tuning the LM’s architecture with task-labeled data, mirroring the process outlined in the first option. The downside of the second option is that data is needed for both the PLM pre-training and the architecture’s fine-tuning. These paradigms are described in detail in the following sections, and an overview is shown in Figure 1.

The paradigm described above requires at least some task-labeled data, but there are some settings where no data is available. A new paradigm for NLP-based modelling has arisen, where an instruction-tuned causal LLM is prompted in natural language to act as an NLP-based model with few or zero examples given (Liu et al., 2023), exploiting its in-context learning ability. This framework is also an option to consider when building NLP-based models.

Data can be restricted in clinical environments for multiple reasons, such as privacy-related issues or the lack of interoperability. These restrictions

Figure 1: Overview of the compatibilities between available data, settings and NLP paradigms that will be described in the paper.



lead to some specific settings, where in one, there is abundant data availability to apply the entire set of NLP modelling paradigms, another where the access to data is incomplete; thus, not all paradigms can be used to develop models, and in the last one no data is available; therefore a specific paradigm should be used. These settings will be carefully described in the following sections and are summarized along with their paradigm compatibility in Figure 1.

Problem A situation arises when there is an asymmetry in data availability or no data is available. In some cases, there is only task-labeled data, only domain-specific unlabeled data, or no data is available at all. Even though multiple paradigms exist for NLP modelling in clinical environments, the compatibility between data availability and the NLP modelling paradigm regarding gains in performance still needs to be explored.

Solution We performed an experimental analysis to measure the performance of solving clinical NLP tasks in Spanish with multiple data availability and NLP modelling paradigm combinations and empirically constructed recommendations for clinical NLP modelling regarding data availability.

1.1 Background

The last paradigm for deep-learning-based NLP was the usage of recurrent neural networks (RNN), which preserved the sequence nature of language in the representation of meaning (Chung et al., 2014; Hochreiter and Schmidhuber, 1997). One drawback of recurrent RNNs is their limited paralleliz-

ability, resulting in prolonged training times. Furthermore, as the sequence lengths grow, there is a tendency for information gathered at distant time steps to vanish due to inherent memory limitations. Nowadays, the Transformer completely ditches the recurrence of the architecture but also preserves word order by learning dependences without regard to their distance in the sentences (Vaswani et al., 2017). With its attention mechanism, this architecture reaches state-of-the-art in multiple NLP tasks such as text classification (Yang et al., 2019), sentiment analysis (Yang et al., 2019), dependency parsing (Mrini et al., 2020), machine translation (Edunov et al., 2018), and named entity recognition (Wang et al., 2021).

1.1.1 Pretrained language models

PLMs are LMs that were trained using self-supervised techniques over large *corpora* of unannotated text to transfer learning from the knowledge gathered in the pre-training to downstream task-specific models (Wang et al., 2022b). Early methods for PLMs consisted of static word embeddings, which were distributed word representations learned using algorithms such as Word2Vec (Mikolov et al., 2013) or GloVe (Pennington et al., 2014), and these embeddings were standard initialization parameters for deep learning architectures to solve NLP tasks. There has been a shift towards dynamic or context-aware word embeddings, which solves the problem of static word embeddings that do not consider word polysemy. These context-aware word embeddings were initially composed using RNNs (Dai and Le, 2015) such as in ELMo (Peters et al., 2018), but currently, they are based on the Transformer architecture and use web-scale unannotated text to be trained. The *de facto* standard for pre-trained Transformer-based context-aware models is BERT (Devlin et al., 2019) and BERT-alike models such as RoBERTa (He et al., 2021) and DeBERTa (He et al., 2021). This language model learns bidirectional contexts conditioning on both left and right contexts in deep stacked layers. Using BERT as a base architecture, domain-specific models have arisen, such as roberta-base-bne (Fandiño et al., 2022), a RoBERTa-based PLM for the Spanish language, PubMedBERT (Gu et al., 2021), a PLM for the biomedical domain and LEGAL-BERT (Chalkidis et al., 2020), a PLM for the legal domain, among others.

165	1.1.2 Large language models	
166	LLMs are PLMs with a significantly larger model	
167	size scale (Zhao et al., 2023); for example, the	
168	PLM BERT has a model size of 0.3×10^9 param-	
169	eters and the LLM GPT-3 (Brown et al., 2020), has	
170	175×10^9 parameters. It has been found that scaling	
171	PLMs improves the performance of the models on	
172	downstream tasks (Kaplan et al., 2020); although	
173	this is true, some other surprising and more im-	
174	portant behaviours in solving a series of complex	
175	tasks appear at LLM scales and were called <i>emer-</i>	
176	<i>gent abilities</i> . Emergent abilities are aptitudes not	
177	present in small models but arise in LLMs (Wei	
178	et al., 2022) and include <i>in-context learning</i> , where	
179	a model can generate expected outputs to natural	
180	language instructions without additional training,	
181	<i>instruction following</i> , where a model fine-tuned	
182	using natural language instructions performs well	
183	on unseen tasks that are also described in the form	
184	of instructions and <i>step-by-step reasoning</i> , where a	
185	model can solve complex problems by instructing	
186	the model involving intermediate reasoning steps	
187	for deriving the final answer. GPT-3, a closed-	
188	source privative LLM, formally introduced the con-	
189	cept of in-context learning, and from there, subse-	
190	quent models have appeared, such as open-source	
191	models Galactica (Taylor et al., 2022), a 120×10^9	
192	parameters model and LLaMA (Touvron et al.,	
193	2023a), a 65×10^9 parameters model. It is worth	
194	noting a significant milestone in LLMs called Chat-	
195	GPT, a closed-source privative assistant-style LLM	
196	that exhibited a superior capacity to communicate	
197	with humans and has been in widespread usage by	
198	laypeople.	
199	1.1.3 Pre-train, fine-tune and predict	
200	paradigm of PLMs	
201	The primary adaptation method for adjusting PLM	
202	to downstream tasks is fine-tuning, where a task-	
203	specific layer is concatenated to the output of	
204	the PLM (Qiu et al., 2020). This method was	
205	proposed in the Universal Language Model Fine-	
206	Tuning (ULMFiT) framework as a transfer learn-	
207	ing technique for domain-specific NLP, achieving	
208	state-of-the-art performances in multiple NLP tasks	
209	(Howard and Ruder, 2018). Even though the fine-	
210	tuning paradigm has been well described for adapt-	
211	ing PLMs, LLMs have significantly higher com-	
212	putational complexity due to their unprecedented	
213	scale. For this reason, some special techniques	
214	have been developed, such as Parameter-Efficient	
215	Fine-Tuning (PEFT), where a small set of param-	
	eters are trained to enable a model to perform the	216
	new task (Ding et al., 2023), showing improve-	217
	ments over in-context learning (Liu et al., 2022).	218
	1.1.4 Pre-train, prompt and predict paradigm	219
	of LLMs	220
	The principal approach for interfacing with LLMs	221
	is through prompting, instructions in natural lan-	222
	guage issued to LLMs to adapt them to new scen-	223
	arios with few or no labelled data (Zhao et al.,	224
	2023) by exploiting the emergent ability of in-	225
	context learning. This new NLP paradigm created a	226
	new field of prompt engineering, where prompting	227
	templates are created to achieve the most effec-	228
	tive performance on downstream tasks (Liu et al.,	229
	2023). There is mixed evidence comparing fine-	230
	tuning vs in-context learning, whereas in some	231
	tasks such as in biomedical information extraction	232
	(Jimenez Gutierrez et al., 2022) or out-of-domain	233
	generalization (Mosbach et al., 2023) fine-tuning	234
	outperforms in-context learning, in other tasks such	235
	as code intelligence (Wang et al., 2022a), in-context	236
	learning outperforms fine-tuning.	237
	1.1.5 Clinical NLP	238
	For clinical NLP, domain-specific models have	239
	been explored in the literature, and their positive	240
	impact on downstream clinical NLP tasks has been	241
	proven (Kalyan and Sangeetha, 2020; Lewis et al.,	242
	2020) even in Spanish (Carrino et al., 2022). There	243
	are public pre-trained Spanish language models for	244
	the clinical domain, including masked LMs, such	245
	as the one we are going to describe in the next sec-	246
	tion and small causal character-level LMs, such as	247
	Clinical-Flair (Rojas et al., 2022), though, in the	248
	large LM category; there are very few, and only for	249
	the English language, such as BioMedLM (Bolton	250
	et al., 2023) and MEDITRON (Chen et al., 2023).	251
	Even though most of the clinical NLP research	252
	has focused on the pre-train, fine-tune, and predict	253
	paradigm, some works have explored the prompt	254
	and predict paradigm through few-shot models	255
	(Sivarajkumar and Wang, 2022), validating that	256
	one can extract clinical information from docu-	257
	ments through prompting (Sivarajkumar et al.,	258
	2023; Agrawal et al., 2022).	259
	2 Data & methods	260
	We intentionally limited data access to evaluate	261
	its impact on the performance of multiple clinical	262
	NLP modelling paradigms and foundation models.	263

264	Each restricted setting was based on a real-world	provider only has access to specific and segmented	313
265	simulated clinical environment.	data sources due to the lack of interoperability.	314
266	2.1 Simulated settings	2.1.3 No data availability	315
267	To mimic clinical settings regarding data availabil-	In this data availability setting, there is no unla-	316
268	ity, we simulated multiple settings with varying	belled unstructured free-text data to continue the	317
269	levels of data availability. We divided the data into	pre-training nor task-specific labelled data to fine-	318
270	two categories: task-specific labelled data, which	tune foundation models. The absence of data can	319
271	can be used to fine-tune models and setting-specific	be attributed to the lack of access to the electronic	320
272	unlabelled data, which can be used to continue the	health record (EHR) database or policies that forbid	321
273	pre-training of the foundation models. The overall	patient data usage to tune machine learning models.	322
274	environment we are located in is a Chilean public	This setting can be seen in a healthcare provider	323
275	health institution analyzing waiting list data, where	using an external EHR service that forbids access	324
276	the explanation of why the patient is waiting is in	to the underlying database, or the provider wants	325
277	the form of free text, and from that dataset, multi-	to solve a new task where data is not yet available.	326
278	ple tasks need to be solved. Multiple reasons can	2.2 Clinical NLP tasks	327
279	restrict data availability; for example, data avail-	To measure the impact of data availability on the	328
280	ability for model training can be restricted due to	performance of clinical NLP modelling, we used	329
281	legal and privacy reasons or because the task trying	multiple clinical NLP tasks, where each is under	330
282	to be solved still does not have sufficient examples	the same environment of the analysis of unstruc-	331
283	due to its recent appearance.	tured waiting list data.	332
284	Unannotated data The unlabelled data we used	2.2.1 Referral prioritization	333
285	to continue the pre-training of the foundation mod-	Different methods exist to prioritize patient selec-	334
286	els was the complete set of reasons for referral con-	tion to process the waiting list more fairly, and	335
287	tained in the Chilean waiting list and is comprised	we modelled the patient prioritization through the	336
288	of 13 365 476 documents, totalling 65 891 568 to-	classification of each referral regarding its state ac-	337
289	kens with a vocabulary size of 513 315 types.	cording to the Chilean Explicit Health Guarantees	338
290	2.1.1 Complete data availability	law (GES in Spanish), which states that specific	339
291	In this data availability setting, unlabelled unstruc-	health problems must be guaranteed to be resolved	340
292	tured free-text data to continue the pre-training and	within a particular time frame. This task requires	341
293	task-specific labelled data are also available to fine-	a binary classification modelling technique. The	342
294	tune foundation models. This setting can be seen	dataset (citation redacted for anonymity) contains	343
295	at a large healthcare provider or at a country-level	1 701 582 examples in the training subset, 485 649	344
296	public health institution such as a ministry of health,	in the test subset and 242 746 in the validation sub-	345
297	where data policies are well established, and pa-	set.	346
298	tients must consent that their data can be used to	2.2.2 Referral speciality classification	347
299	tune machine learning models.	Each referral contained on the waiting list corre-	348
300	2.1.2 Incomplete data availability	sponds to a specific medical speciality. This task	349
301	In this data availability setting, only task-specific la-	involves the prediction of the corresponding medi-	350
302	belled data is available to fine-tune foundation mod-	cal speciality given the free-text description of the	351
303	els. The lack of unlabelled unstructured free-text	reason for referral contained on the waiting list	352
304	data to continue the pre-training may be attributed	record. This task requires a multilabel modelling	353
305	to the fact that the provided is only acquiring data	technique with a label space size of 48 classes. The	354
306	for the specific task and does not have access to	dataset contains 3 401 173 examples in the training	355
307	close-to-the-environment unlabelled text data or ac-	subset, 971 764 in the test subset and 485 882 in	356
308	cording to data policies, the provider cannot merge	the validation subset.	357
309	patient data from a different source, other than the	2.2.3 Clinical named entity recognition	358
310	source of the task data. This setting can be seen at	Clinical named entity recognition is a subtype of	359
311	a medium-sized healthcare provider where the data	named entity recognition in which entities of clin-	360
312	warehousing methods are not implemented or the		

361	ical interest are extracted from unstructured free-	2.3.4 Llama 2	409
362	text sources. This dataset (citation redacted for	Llama 2 is a causal auto-regressive language model	410
363	anonymity) is annotated with eleven different clin-	that uses an optimized transformer architecture	411
364	ical entity classes and was modelled as a token	trained on a <i>corpus</i> of publicly available online data	412
365	classification problem, where each of the tokens	comprised of two trillion tokens (Touvron et al.,	413
366	of the referrals is classified into one of the eleven	2023b). This model is the largest we tested but	414
367	clinical entity classes. The dataset contains 7987	is not domain-adapted in any way, and this is the	415
368	documents in the training subset, 987 in the test	model we used for in-context learning prediction.	416
369	subset and 887 in the validation subset.		
370	2.3 Foundation models	2.4 Modelling paradigms	417
371	We used multiple foundation models as a basis to	We utilized various NLP modelling paradigms to	418
372	solve the clinical NLP tasks. The attributes used	tackle each clinical NLP task, experimenting with	419
373	to select the foundation models were the language,	multiple paradigms for some foundational mod-	420
374	domain and modelling technique.	els based on their compatibility. Also, we note	421
375		the compatibility of each paradigm with each data	422
376	2.3.1 XLM-RoBERTa	availability setting.	423
377	A multilingual version of XLM-RoBERTa masked	2.4.1 Continue pre-training, fine-tune and	424
378	language model, pre-trained using a self-supervised	predict	425
379	technique on a <i>corpus</i> of 2.5TB of filtered Com-	This modelling paradigm is the most data-intensive,	426
380	monCrawl raw text data containing one hundred	where we start with an already pre-trained LM	427
381	languages (Conneau et al., 2019). This model is	checkpoint and continue the pre-training for five	428
382	the broadest of all of our selected foundation LMs.	epochs with the closer-to-the-environment unan-	429
383	This model should be viewed as a baseline where	notated data described in 2.1. Then, with the now	430
384	no model is available for the language or the do-	environment-adapted LM, we perform a fine-tuning	431
385	main.	for five epochs to solve each clinical NLP task. We	432
386	2.3.2 Spanish RoBERTa	continued the pre-training of all the masked LMs	433
387	A Spanish language version of RoBERTa masked	(XLM-RoBERTa, Spanish RoBERTa and Spanish	434
388	language model, pre-trained on a <i>corpus</i> of 570GB	biomedical and clinical RoBERTa) with no modifi-	435
389	of clean and deduplicated text, compiled from the	cation to the original vocabulary and using model-	436
390	web crawlings performed by the National Library	default hyperparameters. This paradigm is com-	437
391	of Spain (Biblioteca Nacional de España) from	patible only with the Complete data availability	438
392	2009 to 2019 (Fandiño et al., 2022). This model	setting.	439
393	is only compatible with the language in which the	2.4.2 Fine-tune and predict	440
394	clinical NLP tasks are and is a type of model (re-	In this paradigm, we started with each of the	441
395	garding language) that should be used when no	off-the-shelf masked foundation models (XLM-	442
396	domain-specific model is available.	RoBERTa, Spanish RoBERTa and Spanish biomed-	443
397		ical and clinical RoBERTa) and performed fine-	444
398	2.3.3 Spanish biomedical and clinical	tuning for each of the clinical NLP tasks. We fine-	445
399	RoBERTa	tuned each task using the default model hyperpa-	446
400	A Spanish language biomedical and clinical version	rameters and trained for five epochs. This paradigm	447
401	of RoBERTa masked language model, pre-trained	is compatible with both Complete data availability	448
402	on a <i>corpus</i> of several biomedical <i>corpora</i> in Span-	and Incomplete data availability settings.	449
403	ish, collected from publicly available <i>corpora</i> and	2.4.3 Prompt and predict	450
404	crawlers, and a real-world clinical <i>corpus</i> . The en-	In this paradigm, we exploited LLMs' in-context	451
405	tire <i>corpus</i> was comprised of more than 1B tokens	learning emergent ability through zero-shot and	452
406	(Carrino et al., 2021). This model is the closest to	few-shot techniques. We prompted the LLM	453
407	the domain model we used to solve the tasks, com-	(Llama 2) to solve each task and parsed its an-	454
408	patible with both language and domain; this should	swer accordingly. For the few-shot technique, we	455
	be the best-suited model to solve a domain-specific	randomly sampled five examples of the training	456
	task.		

Model	Prioritization	Specialty	CNER
xlm-roberta			
Off-the-shelf	88.85 %	51.71 %	11.09 %
Environment-pre-trained	89.03 % (+0.18)	52.36 % (+0.65)	13.85 % (+2.76)
roberta-bne			
Off-the-shelf	88.58 %	52.50 %	22.59 %
Environment-pre-trained	88.80 % (+0.22)	51.65 % (-0.85)	23.29 % (+0.70)
roberta-biomedical-clinical			
Off-the-shelf	88.80 %	53.79 %	34.46 %
Environment-pre-trained	88.85 % (+0.05)	53.85 % (+0.06)	37.25 % (+2.79)
Llama 2			
Zero-shot	6.49 %	31.41 %	5.31 %
Few-shot	56.70 % (+50.21)	31.91 % (+0.50)	15.44 % (+10.13)

Table 1: Results (macro F_1 score) for each clinical NLP task and each fine-tuned model with and without environment continuation of pre-training.

subset of each clinical NLP task. This task is compatible with [Complete data availability](#), [Incomplete data availability](#) and [No data availability](#) settings. The prompt templates used to solve each task are available in the appendix [A](#).

2.5 Increasing training data size and its impact on model performance

To better understand the direct impact of the number of training examples, we performed a test in which we truncated the training subset in increasing steps and measured the performance of the fine-tuned model on the complete test subset. We applied this experiment to all settings and masked LMs.

3 Results

The results for each modelling paradigm solving each clinical NLP task are presented in Table 1.

The [Referral prioritization](#) task was where the model performed the best due to its straightforward binary nature. The models could identify the prioritized health problems mentioned in the training subset and generalize the knowledge correctly in the test subset. On the other hand, the [Clinical named entity recognition](#) task was the most complex of the three tasks, and the models struggled the most to solve it. The performance of the models was directly related to the intrinsic complexity of the clinical NLP task.

Regarding the foundation language model used to solve the clinical NLP tasks, the [Spanish biomedical and clinical RoBERTa](#) model was the best performant; this model is the closest to the domain

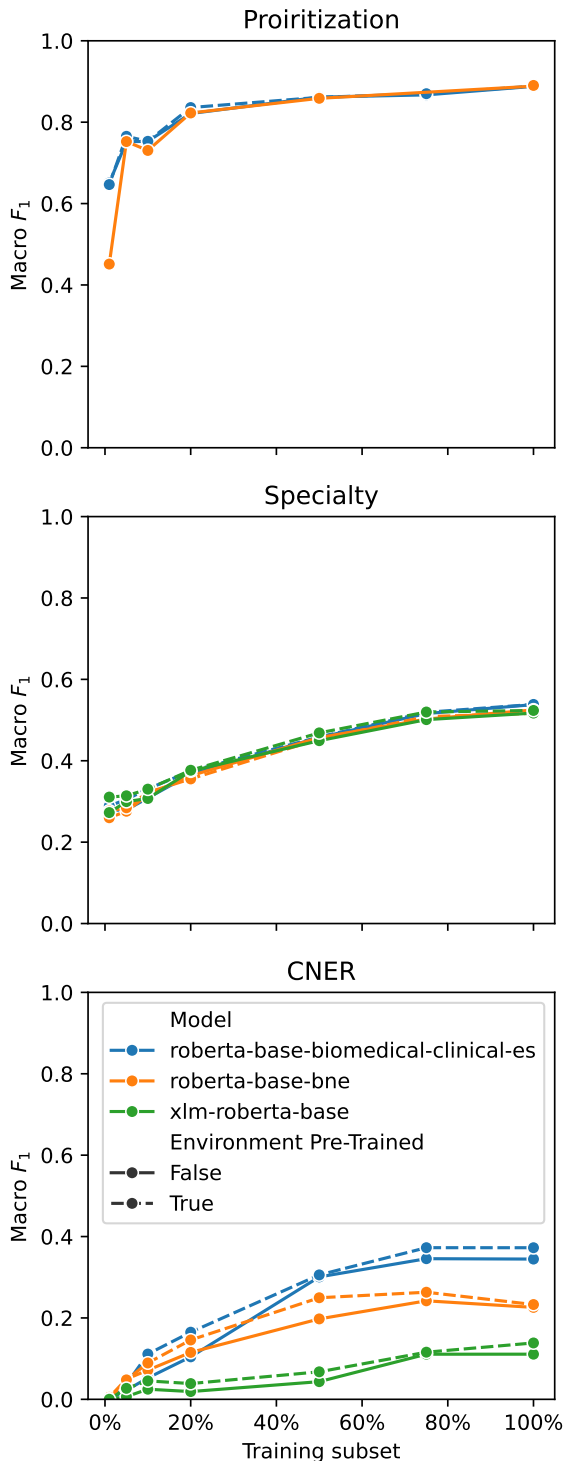
of the clinical NLP tasks and therefore was able to transfer learning from its pre-training on a large clinical corpus. The worst-performing model was [XLM-RoBERTa](#), which is less close to the domain foundation model; therefore, its ability to use prior knowledge to model the tasks was lacking. The closeness to the domain between the foundation model and the task is correlated with model performance on downstream tasks. Before, we only compared models that solved the tasks using the same paradigm, and overall, the worst performant model was [Llama 2](#); however, this model was used by exploiting a different paradigm.

Environment adaptation using unlabelled unstructured free-text data improved model performance. However, the improvements are marginal, considering the computational resources used to continue the pre-training of the foundation models.

The [Continue pre-training, fine-tune and predict](#) paradigm achieved the best result in solving all the clinical NLP tasks. However, we do not recommend using it as a paradigm for clinical NLP modelling, given its high resource usage for training and its overall low gain in performance. On the other hand, the [Prompt and predict](#) paradigm was the worst performant paradigm of all three, but it is worth noting that in some specific cases, its performance was better than the other paradigms. Also, the access to few-shot examples drastically improves in-context learning performance. We recommend using this paradigm in settings with minimal access to training data.

The experiment’s results on the impact of training data volume on the performance of downstream

Figure 2: Performance (macro F_1 score) by training subset for each clinical NLP task and each fine-tuned model with and without environment continuation of pre-training



tasks are presented in Figure 2.

All the models display a performance saturation even before attaining complete training data. This phenomenon is best noticed in the Referral prioritization clinical NLP task, where minimal access

to training data can result in almost peak performance. This behaviour further indicates that the task has a relatively low complexity. The Referral speciality classification task exhibits a more nuanced performance saturation phenomenon than other tasks. The correlation between training data availability and performance is nearly linear, indicating that access to training data is crucial for specific complex tasks.

4 Conclusion

Our study investigated the impact of data availability on the performance of clinical NLP modelling in simulated settings with varying levels of access to task-specific labelled data and unannotated environment-specific text. We explored different paradigms, including Continue pre-training, fine-tune and predict, Fine-tune and predict, as well as Prompt and predict with few-shot learning.

Our findings indicate that choosing foundation models, especially those closer to the target domain impacts model performance. The Spanish biomedical and clinical RoBERTa model, tailored to the clinical domain, outperformed other models in our experiments. While continuing pre-training with environment-specific data improved model performance, the gains were marginal compared to the computational resources required. The fine-tuning paradigm without additional pre-training proved practical, particularly in settings with limited access to unlabelled data.

In-context learning, using the prompt and predict paradigm, demonstrated its viability, especially in scenarios where there is no labelled data available. The creation of few-shot examples significantly improved performance, highlighting the potential of this approach in data-scarce environments.

Our study also revealed a saturation point in performance concerning the amount of training data available. In some instances, minimal data access can still lead to relatively high performance, particularly for less complex tasks.

The choice of foundation models, the utilization of available data, and the selection of appropriate modelling paradigms are crucial considerations in clinical NLP tasks. While pre-training and fine-tuning with domain-specific data remain effective, in-context learning with few-shot examples offers a viable solution in settings where labelled data is unavailable.

577	5 Recommendations		
578	Based on our comprehensive analysis, we provide		
579	recommendations for practitioners engaged in clinical		
580	NLP modelling:		
581	Model selection When selecting foundation models,		
582	prioritize those that align closely with the		
583	target domain. Our results emphasize the significance		
584	of domain specificity in achieving		
585	optimal performance.		
586	Data utilization In settings with ample access		
587	to task-specific labelled data and unannotated		
588	domain-specific text, the Continue pre-training ,		
589	fine-tune and predict paradigm may be considered.		
590	However, given the resource-intensive nature of this		
591	approach, practitioners may opt for the Fine-tune and predict		
592	paradigm, especially when computational resources		
593	are constrained.		
594			
595	If no data is available In scenarios with no access		
596	to labelled data, the Prompt and predict		
597	paradigm, particularly with few-shot learning,		
598	emerges as a practical and effective solution.		
599	This approach allows models to leverage general		
600	knowledge and adapt to new tasks with minimal		
601	labelled examples.		
602	Consideration of task complexity Recognize		
603	the inherent complexity of the clinical NLP		
604	task at hand. Tasks with lower complexity		
605	may achieve near-optimal performance		
606	even with minimal access to training data,		
607	highlighting the importance of task-specific		
608	considerations.		
609	Continuous investigation Clinical NLP is		
610	dynamic, and advancements in pre-trained		
611	foundation LMs and novel paradigms are frequent.		
612	Continuously exploring emerging techniques		
613	and adapting to the evolving landscape is		
614	essential for staying at the forefront of effective		
615	healthcare information extraction.		
616	By incorporating these recommendations, practitioners		
617	can make informed decisions based on the specific		
618	characteristics and constraints of their clinical		
619	NLP projects, ultimately enhancing the efficiency		
620	and efficacy of their models in real-world		
621	healthcare applications.		
	6 Limitations		622
	We attempted to use settings that can be easily		623
	understood in real-world scenarios, but we may		624
	have unintentional biases based on our experiences		625
	in local environments. Our choice of foundation		626
	Language Models (LMs) for each category (multilingual,		627
	language-specific, and domain-specific) may require		628
	a different categorization in order to provide		629
	representative examples of LMs.		630
	7 Ethics statement		631
	We obtained all the data we used through a		632
	transparency law that requires public health providers		633
	to make the data available to the public. This means		634
	that anyone can access the same data that we used,		635
	provided they follow the same process that we did.		636
	The data we used is public, and we have cited the		637
	source papers where each dataset was officially		638
	released to the public.		639
	References		640
	Monica Agrawal, Stefan Hegselmann, Hunter Lang,		641
	Yoon Kim, and David Sontag. 2022. Large language models are few-shot clinical information extractors .		642
	In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages		643
	1998–2022, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.		644
			645
			646
			647
	Elliot Bolton, David Hall, Michihiro Yasunaga, Tony		648
	Lee, Chris Manning, and Percy Liang. 2023. BioMedLM. https://crfm.stanford.edu/2022/12/15/biomedlm.html . [Accessed 14-12-2023].		649
			650
			651
	Tom Brown, Benjamin Mann, Nick Ryder, Melanie		652
	Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind		653
	Neelakantan, Pranav Shyam, Girish Sastry, Amanda		654
	Askell, Sandhini Agarwal, Ariel Herbert-Voss,		655
	Gretchen Krueger, Tom Henighan, Rewon Child,		656
	Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens		657
	Winter, Chris Hesse, Mark Chen, Eric Sigler, Matusz		658
	Litwin, Scott Gray, Benjamin Chess, Jack		659
	Clark, Christopher Berner, Sam McCandlish, Alec		660
	Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 1877–1901. Curran Associates, Inc.		661
			662
			663
			664
			665
	Casimiro Pio Carrino, Jordi Armengol-Estapé, Asier		666
	Gutiérrez-Fandiño, Joan Llop-Palao, Marc Pàmies,		667
	Aitor Gonzalez-Agirre, and Marta Villegas. 2021. Biomedical and clinical language models for spanish: On the benefits of domain-specific pretraining in a mid-resource scenario .		668
			669
			670
			671
	Casimiro Pio Carrino, Joan Llop, Marc Pàmies, Asier		672
	Gutiérrez-Fandiño, Jordi Armengol-Estapé, Joaquín		673

674	Silveira-Ocampo, Alfonso Valencia, Aitor Gonzalez-Agirre, and Marta Villegas. 2022. Pretrained biomedical language models for clinical nlp in spanish . In <i>Proceedings of the 21st Workshop on Biomedical Language Processing</i> . Association for Computational Linguistics.	732
675		733
676		734
677		735
678		736
679		737
680	Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020. LEGAL-BERT: The muppets straight out of law school . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2898–2904, Online. Association for Computational Linguistics.	738
681		739
682		740
683		741
684		742
685		743
686		
687	Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, Alexandre Sallinen, Alireza Sakhaeirad, Vinitra Swamy, Igor Krawczuk, Deniz Bayazit, Axel Marmet, Syrielle Montariol, Mary-Anne Hartley, Martin Jaggi, and Antoine Bosselut. 2023. Meditron-70b: Scaling medical pre-training for large language models .	744
688		745
689		746
690		747
691		748
692		749
693		
694		750
695		751
696		752
697		753
698		
699		754
700		755
701		756
702		
703		757
704		758
705		759
706		760
707		761
708		762
709		
710		763
711		764
712		765
713		766
714		767
715		768
716		769
717		
718		770
719		771
720		772
721		773
722		
723		774
724		775
725		776
726		777
727		778
728		
729		779
730		780
731		781
		782
		783
		784
		785
		786

787	2022. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning . In <i>Advances in Neural Information Processing Systems</i> .	Sonish Sivarajkumar and Yanshan Wang. 2022. Health-Prompt: A zero-shot learning paradigm for clinical natural language processing . <i>AMIA Annu. Symp. Proc.</i> , 2022:972–981.	842
788			843
789			844
790	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing . <i>ACM Computing Surveys</i> , 55(9):1–35.	Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. Galactica: A large language model for science .	846
791			847
792			848
793			849
794			
795	Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality . In <i>Advances in Neural Information Processing Systems</i> , volume 26. Curran Associates, Inc.	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models .	850
796			851
797			852
798			853
799			854
800	Marius Mosbach, Tiago Pimentel, Shauli Ravfogel, Dietrich Klakow, and Yanai Elazar. 2023. Few-shot fine-tuning vs. in-context learning: A fair comparison and evaluation .	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esibou, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models .	855
801			856
802			857
803			858
804	Khalil Mrini, Franck Deroncourt, Quan Hung Tran, Trung Bui, Walter Chang, and Ndapa Nakashole. 2020. Rethinking self-attention: Towards interpretability in neural parsing . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 731–742, Online. Association for Computational Linguistics.		859
805			860
806			861
807			862
808			863
809			864
810			865
811	Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation . In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.		866
812			867
813			868
814			869
815			870
816			871
817	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.		872
818			873
819			874
820			875
821			876
822			877
823			878
824			
825			
826	XiPeng Qiu, TianXiang Sun, YiGe Xu, YunFan Shao, Ning Dai, and XuanJing Huang. 2020. Pre-trained models for natural language processing: A survey . <i>Science China Technological Sciences</i> , 63(10):1872–1897.	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need . In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	879
827			880
828			881
829			882
830			883
831	Matías Rojas, Jocelyn Dunstan, and Fabián Villena. 2022. Clinical flair: A pre-trained language model for Spanish clinical natural language processing . In <i>Proceedings of the 4th Clinical Natural Language Processing Workshop</i> , pages 87–92, Seattle, WA. Association for Computational Linguistics.	Chaozheng Wang, Yuanhang Yang, Cuiyun Gao, Yun Peng, Hongyu Zhang, and Michael R. Lyu. 2022a. No more fine-tuning? an experimental evaluation of prompt tuning in code intelligence . In <i>Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2022</i> , page 382–394, New York, NY, USA. Association for Computing Machinery.	884
832			885
833			886
834			887
835			888
836			889
837	Sonish Sivarajkumar, Mark Kelley, Alyssa Samolyk-Mazzanti, Shyam Visweswaran, and Yanshan Wang. 2023. An empirical evaluation of prompting strategies for large language models in zero-shot clinical natural language processing .	Haifeng Wang, Jiwei Li, Hua Wu, Eduard Hovy, and Yu Sun. 2022b. Pre-trained language models and their applications . <i>Engineering</i> .	890
838			891
839			892
840			893
841			894
		Xinyu Wang, Yong Jiang, Nguyen Bach, Tao Wang, Zhongqiang Huang, Fei Huang, and Kewei Tu. 2021. Automated concatenation of embeddings for structured prediction . In <i>Proceedings of the 59th Annual</i>	895
			896
			897
			898
			899

900 *Meeting of the Association for Computational Lin-*
901 *guistics and the 11th International Joint Conference*
902 *on Natural Language Processing (Volume 1: Long*
903 *Papers)*, pages 2643–2660, Online. Association for
904 Computational Linguistics.

905 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,
906 Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
907 Maarten Bosma, Denny Zhou, Donald Metzler, Ed H.
908 Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy
909 Liang, Jeff Dean, and William Fedus. 2022. [Emer-](#)
910 [gent abilities of large language models](#). *Transactions*
911 *on Machine Learning Research*. Survey Certifica-
912 tion.

913 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Car-
914 bonell, Ruslan Salakhutdinov, and Quoc V. Le. 2019.
915 *XLNet: Generalized Autoregressive Pretraining for*
916 *Language Understanding*. Curran Associates Inc.,
917 Red Hook, NY, USA.

918 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,
919 Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen
920 Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen
921 Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang,
922 Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu,
923 Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023.
924 [A survey of large language models](#). *arXiv preprint*
925 *arXiv:2303.18223*.

926 Zhe Zheng, Xin-Zheng Lu, Ke-Yin Chen, Yu-Cheng
927 Zhou, and Jia-Rui Lin. 2022. [Pretrained domain-](#)
928 [specific language model for natural language process-](#)
929 [ing tasks in the AEC domain](#). *Computers in Industry*,
930 142:103733.

931 **A Prompt templates used for in-context** 932 **learning**

933 We describe the prompt templates we used to solve
934 each clinical NLP task using in-context learning.

935 **A.1 Referral prioritization**

936 **System prompt template** En Chile, las
937 garantías explícitas de salud establecen
938 prioridad para un conjunto de problemas
939 de salud. Debes responder en español sólo
940 la palabra "Verdadero" si la enfermedad
941 que te entregue pertenece a uno de los
942 80 problemas de salud y sólo la palabra
943 "Falso" si la enfermedad no pertenece
944 al conjunto de problemas. Los problemas
945 de salud son: "Accidente Cerebrovascular
946 Isquémico en personas de 15 años y más",
947 "Alivio del dolor y cuidados paliativos
948 por cáncer avanzado ", "Analgesia
949 del Parto", "Artritis Reumatoidea",
950 "Artritis idiopática juvenil", "Asma
951 Bronquial moderada y grave en personas
952 menores de 15 años", "Asma bronquial en

953 personas de 15 años y más", "Cardiopatías
954 congénitas operables en menores de 15
955 años", "Colecistectomía preventiva del
956 cáncer de vesícula en personas de 35 a 49
957 años", "Consumo Perjudicial o Dependencia
958 de riesgo bajo a moderado de alcohol
959 y drogas en personas menores de 20
960 años", "Cáncer Cervicouterino", "Cáncer
961 Colorectal en personas de 15 años y más",
962 "Cáncer Vesical en personas de 15 años
963 y más", "Cáncer de Ovario Epitelial",
964 "Cáncer de mama en personas de 15 años
965 y más", "Cáncer de próstata en personas
966 de 15 años y más", "Cáncer de testículo
967 en personas de 15 años y más", "Cáncer
968 en personas menores de 15 años", "Cáncer
969 gástrico", "Depresión en personas de 15
970 años y más", "Desprendimiento de retina
971 regmatógeno no traumático", "Diabetes
972 Mellitus Tipo 1", "Diabetes Mellitus
973 Tipo 2", "Displasia broncopulmonar
974 del prematuro", "Displasia luxante
975 de caderas", "Disrafias espinales",
976 "Endoprótesis total de cadera en
977 personas de 65 años y más con artrosis de
978 cadera con limitación funcional severa",
979 "Enfermedad Pulmonar Obstructiva
980 Crónica de Tratamiento Ambulatorio",
981 "Enfermedad Renal Crónica Etapa 4 y 5",
982 "Enfermedad de Parkinson", "Epilepsia
983 no refractaria en personas de 15 años
984 y más", "Epilepsia no refractaria en
985 personas desde 1 año y menores de 15 años",
986 "Esclerosis múltiple remitente recurrente
987 ", "Esquizofrenia", "Estrabismo en
988 personas menores de 9 años", "Fibrosis
989 Quística", "Fisura labiopalatina", "Gran
990 Quemado", "Hemofilia", "Hemorragia
991 Subaracnoidea secundaria a Ruptura de
992 Aneurismas Cerebrales", "Hepatitis C",
993 "Hepatitis crónica por Virus Hepatitis
994 B", "Hipertensión arterial primaria o
995 esencial en personas de 15 años y más",
996 "Hipoacusia Bilateral en personas de 65
997 años y más que requieren uso de audífono",
998 "Hipoacusia neurosensorial bilateral del
999 prematuro", "Hipotiroidismo en personas
1000 de 15 años y más", "Infarto agudo del
1001 miocardio", "Infección respiratoria
1002 aguda (IRA) de manejo ambulatorio en
1003 personas menores de 5 años", "Leucemia
1004 en personas de 15 años y más", "Linfomas

1005	en personas de 15 años y más", "Lupus	de salud?.	1056
1006	Eritematoso Sistémico", "Neumonía		
1007	adquirida en la comunidad de manejo	A.2 Referral speciality classification	1057
1008	ambulatorio en personas de 65 años	System prompt template Eres un asistente	1058
1009	y más", "Osteosarcoma en personas	serio que sólo da respuestas precisas	1059
1010	de 15 años y más", "Politraumatizado	y concisas que recibirá diagnósticos en	1060
1011	Grave", "Prevención de Parto Prematuro",	Español y deberás sólo responder con el	1061
1012	"Prevención secundaria enfermedad renal	nombre de la especialidad en Español a	1062
1013	crónica terminal", "Retinopatía del	la cual debe enviarse el diagnóstico.	1063
1014	prematuro", "Retinopatía diabética",	Las especialidades disponibles son:	1064
1015	"Salud Oral Integral del adulto de	TRASTORNOS TEMPOROMANDIBULARES Y DOLOR	1065
1016	60 años", "Salud oral integral de la	OROFACIAL, REHABILITACION: PROTESIS	1066
1017	embarazada", "Salud oral integral	FIJA, NUTRICION, GENETICA, ODONTOLOGIA	1067
1018	para niños y niñas de 6 años",	INDIFERENCIADO, CIRUGIA TORAX,	1068
1019	"Síndrome de Dificultad Respiratoria	CIRUGIA INFANTIL, MEDICINA FAMILIAR,	1069
1020	en el recién nacido", "Síndrome de la	NEUROLOGIA, ONCOLOGIA, OBSTETRICIA,	1070
1021	inmunodeficiencia adquirida VIH/SIDA",	CIRUGIA ADULTO, DERMATOLOGIA,	1071
1022	"Trastorno Bipolar en personas de 15	GERIATRIA, OTORRINOLARINGOLOGIA,	1072
1023	años y más", "Trastornos de generación	BRONCOPULMONAR, MEDICINA INTERNA,	1073
1024	del impulso y conducción en personas	PERIODONCIA, CARDIOLOGIA, OFTALMOLOGIA,	1074
1025	de 15 años y más, que requieren	REHABILITACION: PROTESIS REMOVIBLE,	1075
1026	Marcapaso", "Tratamiento Médico en	ENDOCRINOLOGIA, PEDIATRIA, REUMATOLOGIA,	1076
1027	personas de 55 años y más con Artrosis	CIRUGIA PLASTICA, ORTODONCIA, CIRUGIA	1077
1028	de Cadera y/o Rodilla, leve o moderada",	DE MAMAS, CIRUGIA PROCTOLOGICA,	1078
1029	"Tratamiento Quirúrgico de Hernia del	GASTROENTEROLOGIA, HEMATOLOGIA,	1079
1030	Núcleo Pulposo Lumbar", "Tratamiento	UROLOGIA, ANESTESIOLOGIA, ENFERMEDADES	1080
1031	Quirúrgico de lesiones crónicas de la	DE TRANSMISION SEXUAL, OPERATORIA,	1081
1032	válvula aórtica en personas de 15 años	NEONATOLOGIA, NEUROCIRUGIA, CIRUGIA	1082
1033	y más", "Tratamiento Quirúrgico de	VASCULAR PERIFERICA, GINECOLOGIA, CIRUGIA	1083
1034	lesiones crónicas de las válvulas mitral	BUCAL, CIRUGIA MAXILO FACIAL, CIRUGIA	1084
1035	y tricúspide en personas de 15 años	ABDOMINAL, CARDIOCIRUGIA, PSIQUIATRIA,	1085
1036	y más", "Tratamiento de Erradicación	INFECTOLOGIA, TRAUMATOLOGIA, ENDODONCIA,	1086
1037	del Helicobacter Pylori", "Tratamiento	MEDICINA FISICA Y REHABILITACION,	1087
1038	de Hipoacusia moderada en personas	NEFROLOGIA.	1088
1039	menores de 4 años", "Tratamiento de la		
1040	hiperplasia benigna de la próstata en	User prompt template ¿A qué especialidad	1089
1041	personas sintomáticas", "Tratamiento	debo enviar el diagnóstico "<x>"?.	1090
1042	quirúrgico de cataratas", "Tratamiento		
1043	quirúrgico de escoliosis en personas	A.3 Clinical named entity recognition	1091
1044	menores de 25 años", "Trauma Ocular	System prompt template Eres reconecedor	1092
1045	Grave", "Traumatismo Cráneo Encefálico	de entidades nombradas que solo debe	1093
1046	moderado o grave", "Tumores Primarios	detectar las entidades en la siguiente	1094
1047	del Sistema Nervioso Central en personas	lista: "disease": "alteracion o	1095
1048	de 15 años o más", "Urgencia Odontológica	desviacion del estado fisiologico en una	1096
1049	Ambulatoria", "Vicios de refracción en	o varias partes del cuerpo, por causas	1097
1050	personas de 65 años y más" y "Órtesis	en general conocidas, manifestada por	1098
1051	(o ayudas técnicas) para personas de 65	sintomas y signos caracteristicos, y cuya	1099
1052	años y más"	evolucion es mas o menos previsible",	1100
1053	User prompt template ¿"<x>" pertenece a	- medication: "Medicamentos o drogas	1101
1054	la lista de 80 problemas de salud	empleadas en el tratamiento y o prevención	1102
1055	priorizados por las garantías explícitas	de enfermedades", - abbreviation:	1103
		"Abreviatura", - body_part: "Órgano o	1104

1105 una parte anatómica de una persona", -
1106 family_member: "Miembro de la familia", -
1107 laboratory_or_test_result: "Resultado de
1108 laboratorio o test", - clinical_finding:
1109 "Observaciones, juicios o evaluaciones
1110 que se hacen sobre los pacientes", -
1111 diagnostic_procedure: "Exámenes que
1112 permiten determinar la condición del
1113 individuo ", - laboratory_procedure:
1114 "Exámenes que se realizan en diversas
1115 muestras de pacientes que permiten
1116 diagnosticar enfermedades mediante la
1117 detección de biomarcadores y otros
1118 parámetros", - therapeutic_procedure:
1119 "Actividad o tratamiento que es empleado
1120 para prevenir, reparar, eliminar o
1121 curar la enfermedad del individuo",
1122 Debes responder con el mismo texto
1123 de entrada, pero con las entidades
1124 nombradas anotadas con etiquetas en
1125 la misma línea (<nombre_entidad>lorem
1126 ipsum</nombre_entidad>), donde cada
1127 etiqueta corresponde a un nombre de
1128 entidad, por ejemplo: <entidad>Sed
1129 ut perspiciatis</entidad> unde omnis
1130 iste natus error sit voluptatem
1131 <entidad>accusantium</entidad>.
1132 Las únicas etiquetas disponibles
1133 son: medication, abbreviation,
1134 body_part, family_member,
1135 laboratory_or_test_result,
1136 clinical_finding, diagnostic_procedure,
1137 laboratory_procedure,
1138 therapeutic_procedure, no puedes agregar
1139 más etiquetas de las incluidas en esa
1140 lista. IMPORTANTE: NO DEBES CAMBIAR
1141 EL TEXTO DE ENTRADA, SÓLO AGREGAR LAS
1142 ETIQUETAS.