

Cecilia: El Modelo de Lenguaje Cubano

0.1 Introducción

El desarrollo de modelos de lenguaje específicos para variantes regionales del español es una necesidad apremiante en el campo del Procesamiento de Lenguaje Natural (PLN). Los modelos generalistas, aunque poderosos, presentan limitaciones notables cuando se aplican a lenguas con pocos recursos o a variantes regionales, ya que suelen estar entrenados principalmente con datos de idiomas dominantes y no logran capturar los matices lingüísticos, culturales y contextuales propios de comunidades específicas[9]. En el caso del español cubano, esta carencia es especialmente crítica: no existen modelos de lenguaje que reflejen de manera precisa las particularidades léxicas, sintácticas y pragmáticas de la variante cubana, lo que limita el desarrollo de aplicaciones tecnológicas relevantes para Cuba y su diáspora. La creación de un modelo adaptado a esta realidad no solo contribuiría a cerrar la brecha digital lingüística, sino que también permitiría preservar y potenciar la riqueza cultural del español cubano en el ecosistema digital.

La inexistencia de un modelo de lenguaje entrenado específicamente para el español cubano implica que las expresiones idiomáticas, referencias culturales y fenómenos lingüísticos propios de la isla no son comprendidos ni representados adecuadamente por los modelos actuales. Esta situación afecta negativamente el rendimiento de tareas como el análisis de sentimiento, la generación de texto, la traducción automática y la interacción conversacional en contextos cubanos, donde la identidad lingüística es un componente esencial de la comunicación. Afortunadamente, los avances recientes en técnicas de preentrenamiento continuo permiten aprovechar modelos multilingües robustos —como Salamandra 2B— y especializarlos eficientemente en dominios regionales mediante la exposición a corpus representativos, sin necesidad de modificar la arquitectura ni el tokenizador original. Esto abre la puerta a la creación de modelos pequeños, eficientes y culturalmente adaptados, viables incluso en entornos con recursos computacionales limitados.

En este artículo se presenta **Cecilia 2b**, el primer modelo de lenguaje entrenado específicamente para el español cubano. Cecilia 2b está basado íntegramente en la arquitectura de Salamandra 2B, un modelo multilingüe de 2 mil millones de parámetros, y ha sido adaptado a través de preentrenamiento continuo sobre un corpus cuidadosamente compilado de textos cubanos. El corpus incluye más de 296,000 documentos y cerca de mil millones de tokens

provenientes de prensa nacional, literatura, legislación, enciclopedias y letras de canciones, asegurando una cobertura amplia de registros y contextos socioculturales.

El modelo mantiene la compatibilidad total con el diseño y el vocabulario de Salamandra 2B, garantizando que las mejoras observadas se deban exclusivamente a la especialización en datos cubanos. Cecilia 2b es, por tanto, una herramienta pionera que facilita el desarrollo de aplicaciones de PLN adaptadas a la realidad lingüística y cultural de Cuba, sentando las bases para futuros avances en la inclusión digital y la preservación del patrimonio lingüístico regional.

El artículo se organiza de la siguiente manera: en la Sección 2 se revisa el estado del arte en modelos de lenguaje para variantes regionales y lenguas con pocos recursos, destacando las limitaciones de los modelos generalistas y las ventajas de los modelos pequeños y adaptados. La Sección 3 describe en detalle la arquitectura de Cecilia 2b, el proceso de construcción del corpus cubano y los procedimientos de preentrenamiento continuo empleados. En la Sección 4 se presentan los resultados de la evaluación cuantitativa del modelo en tareas estándar multi-lingües, discutiendo tanto su robustez general como las limitaciones actuales en la evaluación específica del español cubano. Finalmente, en la Sección 5 se exponen las conclusiones y se plantean las líneas futuras de trabajo, enfatizando la necesidad de recursos de evaluación y aplicaciones prácticas que validen el impacto de Cecilia 2b en la comunidad cubana.

1 Estado del Arte

1.1 Limitaciones de los modelos generalistas

Los modelos de lenguaje de gran escala (LLMs) han revolucionado el procesamiento del lenguaje natural, pero presentan limitaciones significativas cuando se aplican a lenguajes con pocos recursos o variantes regionales específicas. Estas limitaciones afectan tanto la representación como la funcionalidad de estos modelos en contextos lingüísticos diversos.

1.1.1 Escasez de datos etiquetados y no etiquetados

Los lenguajes con pocos recursos enfrentan dos limitaciones cruciales: una escasez de datos lingüísticos etiquetados y no etiquetados, y la baja calidad de los datos disponibles, que frecuentemente no son suficientemente representativos de los idiomas y sus contextos socioculturales. Esta carencia de datos de entrenamiento adecuados resulta en modelos que no capturan correctamente las particularidades lingüísticas de estos idiomas, lo que afecta su rendimiento en tareas básicas de procesamiento de lenguaje natural.

Los modelos generalistas, entrenados principalmente con datos de idiomas dominantes como el inglés, muestran un rendimiento significativamente inferior cuando se aplican a lenguajes con recursos limitados. Por ejemplo, investigaciones recientes han demostrado que incluso modelos

avanzados como GPT-4o y Llama 3.1 (405B) tienen un rendimiento inferior en comparación con modelos BERT ajustados específicamente para idiomas como el marathi, con márgenes de precisión de 10.2% y 14.1% respectivamente[9].

1.1.2 Limitaciones en la comprensión contextual y cultural

Los modelos generalistas también presentan limitaciones significativas en la comprensión de contextos culturales específicos. Estos modelos carecen de la capacidad para interpretar adecuadamente expresiones idiomáticas, referencias culturales y matices lingüísticos propios de variantes regionales. Esta deficiencia se manifiesta en una menor precisión en tareas como el análisis de sentimiento, la detección de discurso de odio y la clasificación de textos en idiomas con pocos recursos.

Las evaluaciones existentes para lenguajes con pocos recursos contienen limitaciones que necesitan ser estudiadas más a fondo, ya que los marcos de evaluación actuales no capturan adecuadamente las inconsistencias culturales en los conjuntos de datos[4]. Esto resulta en una disminución del rendimiento cuando se aplican modelos multilingües a contextos lingüísticos específicos.

1.1.3 Barreras tecnológicas y de infraestructura

Los lenguajes con pocos recursos a menudo carecen de la infraestructura digital y las herramientas necesarias (como tokenizadores, analizadores morfológicos) que están fácilmente disponibles para lenguajes con muchos recursos. Esta carencia complica incluso los pasos iniciales del desarrollo de modelos, creando barreras tecnológicas significativas para la inclusión de estos idiomas en el panorama de la IA.

La falta de representación adecuada en los modelos generalistas contribuye a ampliar la brecha digital lingüística, dejando a muchas comunidades lingüísticas sin acceso a tecnologías basadas en IA que podrían beneficiar su desarrollo educativo, cultural y económico.

1.2 Ventajas de modelos pequeños para lenguajes regionales

Los modelos de lenguaje pequeños (SLMs) ofrecen ventajas significativas para el desarrollo de soluciones lingüísticas regionales, presentando alternativas viables a los grandes modelos generalistas para contextos específicos.

1.2.1 Eficiencia computacional y accesibilidad

Los SLMs requieren significativamente menos memoria, almacenamiento y potencia de procesamiento en comparación con los modelos grandes, lo que los hace adecuados para dispositivos con capacidades de hardware limitadas, como smartphones, tablets y dispositivos IoT[14]. Esta característica es particularmente relevante en regiones donde el acceso a infraestructura computacional avanzada es limitado.

La eficiencia energética de los SLMs representa otra ventaja crucial, ya que consumen menos energía que sus contrapartes más grandes. Esto no solo resulta rentable sino que también apoya el desarrollo de IA ambientalmente sostenible al reducir la huella de carbono de los modelos de aprendizaje automático. Para comunidades con recursos limitados, esta eficiencia energética puede ser determinante para la adopción de tecnologías de IA.

1.2.2 Adaptabilidad a contextos lingüísticos específicos

Los modelos pequeños son más fáciles de ajustar para aplicaciones específicas y dominios lingüísticos particulares[7]. Al no requerir grandes cantidades de datos o potencia computacional para su entrenamiento, los desarrolladores pueden adaptarlos para tareas como el procesamiento de variantes dialectales, expresiones regionales y contextos culturales específicos.

Investigaciones recientes han demostrado que los SLMs pueden proporcionar una comprensión del lenguaje de alta calidad con un consumo de recursos significativamente menor, lo que los hace ideales para habilitar el trabajo digital en contextos lingüísticos específicos[3]. Su capacidad para capturar dialectos locales, expresiones idiomáticas y matices culturales con mayor precisión que los modelos generalizados grandes contribuye no solo a la preservación del lenguaje sino también a mejorar la inclusión digital en regiones desatendidas.

1.2.3 Mejoras en privacidad y control local

Los SLMs pueden desplegarse localmente, lo que permite el procesamiento de datos en el dispositivo y reduce la necesidad de enviar información sensible a sistemas basados en la nube. Esta característica ofrece ventajas significativas de privacidad para comunidades lingüísticas que pueden tener preocupaciones sobre la soberanía de sus datos culturales y lingüísticos.

Para las lenguas indígenas y regionales, los SLMs ofrecen una solución rentable y eficiente en términos de recursos al reducir los requisitos computacionales y de datos, mientras mejoran la precisión de salida a través de conjuntos de datos específicos y contextualizados. Este enfoque permite un desarrollo más participativo, donde las comunidades lingüísticas pueden mantener mayor control sobre sus recursos lingüísticos y culturales.

1.3 Estrategias para construir modelos de lenguaje pequeños con pocos recursos

El desarrollo de modelos de lenguaje para idiomas con recursos limitados requiere enfoques innovadores que maximicen la eficiencia y efectividad del entrenamiento. Las siguientes estrategias han demostrado ser particularmente efectivas en este contexto.

1.3.1 Continual pretraining como estrategia principal

El preentrenamiento continuo (continual pretraining) ofrece un camino prometedor para la adaptación de dominio con recursos computacionales limitados[16]. Esta técnica permite que modelos preentrenados existentes sean posteriormente entrenados con datos específicos de dominio, permitiéndoles adquirir conocimiento especializado mientras aprovechan su base de conocimiento existente.

Investigaciones recientes han demostrado mejoras significativas en el rendimiento a través del entrenamiento incremental en 400 millones de tokens, seguido de entrenamiento adicional para alcanzar mil millones de tokens. Los resultados muestran ganancias notables en tareas intensivas en conocimiento (MMLU +8.1%) y comprensión contextual (HellaSwag +7.6%), mientras revelan compensaciones en la especialización de dominio[8].

El preentrenamiento continuo de modelos de lenguaje pequeños en corpus específicos de dominio ha demostrado ser más efectivo que entrenar modelos desde cero. Por ejemplo, en el dominio biomédico, los modelos inicializados con MiniLM y continuamente preentrenados en textos específicos del dominio superaron a los modelos entrenados desde cero con el mismo vocabulario[17].

1.3.2 Estrategias de transferencia y adaptación

El enfoque “Adapt-and-Distill” representa una estrategia efectiva para desarrollar modelos pequeños, rápidos y efectivos para dominios específicos. Este método combina la adaptación de modelos preentrenados generales y la destilación de conocimiento específico del dominio, logrando un mejor rendimiento mientras se reduce significativamente el tamaño y se aumenta la velocidad del modelo[18].

La expansión de vocabulario específico del dominio durante la fase de adaptación y el empleo de la probabilidad de ocurrencia a nivel de corpus para elegir automáticamente el tamaño del vocabulario incremental son técnicas clave en este enfoque. Experimentos en los dominios biomédico e informático han demostrado que esta estrategia logra un mejor rendimiento en tareas específicas del dominio mientras el modelo es 3.3 veces más pequeño y 5.1 veces más rápido que los modelos originales[18].

1.3.3 Técnicas para escenarios de recursos extremadamente limitados

Para lenguajes con recursos extremadamente limitados, el enfoque de “datos pequeños” ha demostrado ser sorprendentemente efectivo. Investigaciones recientes han desafiado la suposición común de que las lenguas con pocos recursos se benefician del entrenamiento conjunto con lenguas de mayores recursos, demostrando que es posible entrenar modelos de lenguaje multilingües competitivos con menos de 1 GB de texto[12].

La combinación de datos sintéticos generados tanto por traducción automática estadística como por modelos de traducción automática neuronal multilingües ha demostrado mejorar el rendimiento para lenguas con pocos recursos debido a la mayor diversidad de los datos sintéticos generados. Esta técnica es particularmente valiosa cuando los datos paralelos bilingües son escasos[13].

El uso de técnicas eficientes en parámetros como LoRA PEFT (Parameter-Efficient Fine-Tuning) minimiza el número de parámetros durante el ajuste fino, ofreciendo eficiencia computacional y manteniendo la robustez del modelo original al ajustar solo algunos de los parámetros[2]. Estudios más amplios han enfatizado que el uso de LoRA en entornos con pocos recursos conlleva una baja sobrecarga computacional.

1.4 Proyectos Regionales de Modelos de Lenguaje

El proyecto SEALD (Southeast Asian Languages in One Network Data) constituye una de las iniciativas más ambiciosas para fortalecer la presencia digital de las lenguas del Sudeste Asiático. Mediante la colaboración entre AI Singapore y Google Research, se recopilieron y curaron grandes volúmenes de datos multilingües, abarcando idiomas como indonesio, malayo, tamil, birmano, filipino, vietnamita, tailandés, lao y jemer[11]. Este esfuerzo permitió el desarrollo de SEA-LION, una familia de modelos de lenguaje preentrenados específicamente para la región, con arquitecturas de 3 a 7 mil millones de parámetros y un vocabulario adaptado a las características lingüísticas del área, mejorando sustancialmente la comprensión y generación de texto en estos idiomas.

AfriBERTa[12] representa un enfoque innovador para lenguas africanas con pocos recursos, desafiando la suposición de que el entrenamiento conjunto con idiomas de alto recurso es siempre beneficioso. Este modelo fue entrenado exclusivamente con menos de 1 GB de texto de 11 lenguas africanas, incluyendo el primer modelo de lenguaje para cuatro de ellas. AfriBERTa demostró, en tareas de reconocimiento de entidades nombradas y clasificación de texto, que un modelo multilingüe focalizado puede superar a alternativas generalistas como mBERT y XLM-R, validando la eficacia de estrategias centradas en corpus pequeños y específicos.

Cecilia es el primer modelo de lenguaje cubano, desarrollado a partir de un proceso de preentrenamiento continuo sobre texto escrito cubano. El modelo fue entrenado con aproximadamente 1,000 millones de tokens provenientes de fuentes representativas: prensa nacional, la enciclopedia Ecured, legislación, literatura, enciclopedias de cubanismos y letras de canciones.

Este enfoque permitió a Cecilia capturar matices lingüísticos y culturales propios del español cubano, facilitando aplicaciones de PLN adaptadas a la realidad local.

Salamandra[5] es un caso paradigmático de éxito en la construcción de modelos multilingües europeos, sirviendo también como base para adaptaciones regionales como Cecilia. La arquitectura de Salamandra abarca variantes de 2, 7 y 40 mil millones de parámetros, todas entrenadas desde cero sobre un corpus multilingüe cuidadosamente curado de 7.8 billones de tokens en 35 idiomas europeos y código de programación. El modelo utiliza precisión bfloat16, embeddings RoPE, activación SwiGLU, normalización RMS, atención flash y una longitud de contexto de hasta 8,192 tokens, con un vocabulario de 256,000 tokens. Entrenado en el supercomputador MareNostrum 5, Salamandra ha demostrado un rendimiento competitivo en benchmarks multilingües y sirve como plataforma robusta para la especialización en variantes lingüísticas regionales.

2 Arquitectura y Entrenamiento de Cecilia

Cecilia 2b está basada íntegramente en la arquitectura de Salamandra 2B (BSC-LT/salamandra-2b), un modelo de lenguaje multilingüe de 2 mil millones de parámetros. Para la adaptación al español cubano, no se realizaron modificaciones ni en la arquitectura original ni en el tokenizador empleado, manteniendo la compatibilidad total con el diseño, el tamaño de vocabulario y las capacidades de representación del modelo base. Esta decisión asegura que las mejoras en desempeño se deban exclusivamente al preentrenamiento continuo sobre datos cubanos, y no a cambios estructurales o de tokenización.

2.1 Descripción del Corpus

El corpus utilizado para el preentrenamiento de Cecilia fue cuidadosamente compilado para capturar la diversidad temática y cultural del español cubano. Incluye textos provenientes de múltiples dominios en 296,311 archivos, abarcando fuentes como Ecured, la Enciclopedia Digital del Audiovisual Cubano, prensa nacional, literatura cubana, legislación, enciclopedias de cubanismos, y letras de canciones populares. Esta composición asegura la representación de registros formales e informales, así como de distintos géneros discursivos y contextos socio-culturales, proporcionando una base sólida para que el modelo aprenda matices lingüísticos y culturales propios de Cuba.

En términos cuantitativos, el corpus contiene un total de 2,631,691,355 caracteres y 384,963,687 palabras distribuidas en 34,505,341 líneas. La longitud promedio por documento es de 8,881 caracteres, mientras que la longitud promedio de las oraciones es de 17 palabras. La densidad léxica, medida como caracteres por palabra, es de 6.8, reflejando la riqueza y variedad del vocabulario presente en el dataset. Estos valores aseguran una cobertura amplia tanto en extensión como en profundidad temática.

Métrica	Valor
Archivos totales	296,311
Caracteres totales	2,631,691,355
Palabras totales	384,963,687
Líneas totales	34,505,341
Longitud promedio doc	8,881 caracteres
Longitud promedio oración	17.0 palabras
Densidad léxica	6.8 car./palabra

Para el proceso de entrenamiento, el corpus fue tokenizado en secuencias de hasta 1,024 tokens utilizando el tokenizador original de Salamandra 2B, sin modificaciones[15]. El dataset tokenizado resultante consta de 1,104,532 muestras, sumando 982,024,795 tokens efectivos (sin padding) y 1,131,040,768 tokens totales (con padding). La longitud promedio de secuencia es de 889.3 tokens, y la razón de padding es eficiente, con un 13.2%. Este procesamiento garantiza que la mayor parte de la capacidad del modelo se utilice en datos reales y representativos.

Métrica	Valor
Muestras totales	1,104,532
Tokens totales (sin padding)	982,024,795
Tokens totales (con padding)	1,131,040,768
Longitud promedio de secuencia (tokens)	889.3
Proporción de padding	13.2%

El corpus fue sometido a un riguroso proceso de limpieza y validación, asegurando la ausencia de problemas de encoding y una alta integridad textual. Esto garantiza que los datos de entrada sean de calidad y representativos, fundamentales para un preentrenamiento efectivo y robusto del modelo.

2.2 Proceso de Entrenamiento

El entrenamiento de Cecilia 2b se realizó mediante preentrenamiento continuo sobre el corpus cubano, manteniendo la arquitectura y el tokenizador originales de Salamandra 2B. A continuación se detallan los hiperparámetros empleados:

Parámetro	Valor
Número de épocas	2
Batch size	4
Gradient accumulation steps	16
Batch efectivo	64

Parámetro	Valor
Learning rate	2e-5
Scheduler	Warmup linear
Proporción de warmup	6%
Optimizador	AdamW
Weight decay	0.01
Betas (1, 2)	0.9, 0.999
Gradient clipping norm	1.0
Precisión	bfloat16

Explicación de los hiperparámetros:

- **Número de épocas:** El corpus completo se procesó dos veces, permitiendo al modelo refinar sus representaciones sobre los datos cubanos.
- **Batch size y acumulación de gradientes:** Se utilizaron lotes pequeños (4 ejemplos por paso), acumulando gradientes durante 16 pasos para lograr un batch efectivo de 64, lo que mejora la estabilidad y eficiencia del entrenamiento en hardware limitado.
- **Learning rate y scheduler:** Se empleó una tasa de aprendizaje inicial baja (2e-5) con un scheduler de tipo “warmup linear”, que incrementa gradualmente la tasa de aprendizaje durante el 6% inicial de los pasos de entrenamiento antes de decaer linealmente, ayudando a evitar inestabilidades al comienzo.
- **Optimizador AdamW:** AdamW es una variante de Adam que desacopla el weight decay, permitiendo un mejor control de la regularización y ayudando a prevenir el sobreajuste.
- **Weight decay y betas:** Weight decay de 0.01 promueve la regularización. Los parámetros 1 y 2 controlan los promedios móviles de los momentos de primer y segundo orden en AdamW.
- **Gradient clipping norm:** Limita la norma de los gradientes a 1.0 para evitar explosiones de gradiente.
- **Precisión bfloat16:** Permite mayor eficiencia en memoria y cómputo, facilitando el entrenamiento de modelos grandes en hardware moderno.

Durante el entrenamiento se emplearon técnicas avanzadas como Fully Sharded Data Parallel (FSDP) y gradient checkpointing para maximizar el uso eficiente de memoria y recursos computacionales, además de validaciones periódicas y monitoreo de métricas para garantizar la robustez y generalización del modelo[10].

3 Evaluación

Para evaluar el impacto del preentrenamiento continuo de Cecilia 2b sobre datos cubanos, se comparó su desempeño con Salamandra 2B en una batería de tareas clásicas de procesamiento

de lenguaje natural en inglés y español. Es importante destacar que estas tareas no están diseñadas para medir la comprensión de la variante cubana del español, sino que sirven para verificar la robustez general del modelo tras la adaptación.

Table 4: Resultados de la evaluación cuantitativa en comparación con Salamandra 2B.

Task	Metric	Salamandra	Cecila	Rel Err
arc_challenge	acc	0.37031	0.38225	3.13%
arc_easy	acc	0.72264	0.73401	1.55%
belebele_en	acc	0.21556	0.24778	13.00%
belebele_es	acc	0.22778	0.24444	6.82%
escola	acc	0.59259	0.55461	-6.41%
openbookqa	acc	0.30000	0.28200	-6.00%
openbookqa_es	acc	0.30800	0.29400	-4.55%
paws_en	acc	0.56100	0.57350	2.18%
paws_es	acc	0.56050	0.55550	-0.89%
piqa	acc	0.73721	0.73667	-0.07%
social_iqa	acc	0.45394	0.44626	-1.69%
teca	acc	0.46481	0.43174	-7.11%
wnli	acc	0.46479	0.42254	-9.09%
wnli_es	acc	0.56338	0.59155	4.76%
xnli_en	acc	0.46225	0.47671	3.03%
xnli_va	acc	0.47505	0.48523	2.10%
xstorycloze_en	acc	0.71145	0.70483	-0.93%
xstorycloze_es	acc	0.65255	0.65189	-0.10%
arc_challenge	acc_norm	0.40700	0.41809	2.65%
arc_easy	acc_norm	0.72559	0.73990	1.93%
belebele_en	acc_norm	0.21556	0.24778	13.00%
belebele_es	acc_norm	0.22778	0.24444	6.82%
openbookqa	acc_norm	0.39600	0.40000	1.00%
openbookqa_es	acc_norm	0.40800	0.40400	-0.98%
piqa	acc_norm	0.74701	0.74701	0.00%
cocoteros_es	bleu	8.46507	6.72269	-20.58%
xlsum_es	bleu	0.80082	0.59723	-25.42%
triviaqa	exact_match	0.37595	0.35432	-5.75%
xquad_es	exact_match	0.37731	0.36050	-4.45%
xquad_es	f1	0.58413	0.56911	-2.57%
cocoteros_es	rouge1	0.33887	0.31209	-7.90%
xlsum_es	rouge1	0.13464	0.08705	-35.35%
Mean Diff				-2.43%

La evaluación de Cecilia 2b se realizó sobre un conjunto de tareas estándar ampliamente utilizadas en la comunidad de PLN, empleando benchmarks reconocidos como SpanishBench, CatalanBench, BasqueBench, GalicianBench y tareas en inglés del LM Evaluation Harness. Estas colecciones agrupan tareas de comprensión lectora (por ejemplo, Belebele, que consiste en responder preguntas de opción múltiple sobre textos breves), inferencia lógica (XNLI, WNLI), razonamiento de sentido común (XStoryCloze), identificación de paráfrasis (PAWS), y preguntas de respuesta abierta y razonamiento científico (OpenBookQA, ARC Challenge y ARC Easy). Para las tareas en español, los datasets utilizados son versiones traducidas profesionalmente o generadas y revisadas por humanos, garantizando alta calidad y relevancia para la evaluación de modelos multilingües.

Además, se incluyen benchmarks de traducción automática (como Flores), tareas de resumen (XLSum), y comprensión y respuesta a preguntas de trivia (TriviaQA, XQuAD), cubriendo así un espectro amplio de habilidades lingüísticas y cognitivas. Estos benchmarks permiten comparar el desempeño general de modelos multilingües y monolingües en tareas de comprensión, inferencia, generación y traducción en diferentes idiomas, aunque, como se señaló previamente, no están diseñados para evaluar competencias específicas en variantes regionales como el español cubano.

En promedio, la reducción relativa de desempeño es de apenas 2.4% respecto al modelo base, una diferencia no significativa considerando la magnitud del cambio en los datos de entrenamiento y la especialización lograda[6]. Esto indica que el modelo mantiene su capacidad general para tareas estándar, a pesar de haber sido adaptado a un dominio lingüístico y cultural específico.

Sin embargo, aún no se dispone de benchmarks ni corpus de instrucciones diseñados para evaluar específicamente la comprensión y generación en español cubano. La creación de estos recursos será esencial para medir el verdadero valor añadido de Cecilia en aplicaciones donde el dominio cultural y lingüístico local es crítico.

3.1 Discusión

La principal limitación de Cecilia 2b radica en que su evaluación actual se ha realizado exclusivamente sobre tareas generales y benchmarks estándar de procesamiento de lenguaje natural, como comprensión lectora, inferencia lógica y razonamiento de sentido común, que no están diseñados para medir la comprensión ni la generación en la variante cubana del español. Esto significa que, aunque el modelo mantiene un desempeño robusto en tareas multilingües generales con solo una reducción promedio del 2.4% respecto al modelo base, aún no es posible cuantificar su ventaja específica en aplicaciones donde los matices culturales y lingüísticos cubanos sean críticos.

Además, el corpus utilizado, si bien diverso y representativo, podría inducir sesgos hacia los dominios más representados, como prensa y enciclopedias, en detrimento de registros menos frecuentes o más formales, lo que podría limitar la cobertura de ciertos contextos y estilos

lingüísticos[1]. Esta situación resalta la necesidad de continuar ampliando y balanceando el corpus, incorporando textos de géneros subrepresentados y registros formales, así como de desarrollar recursos de evaluación específicos para el español cubano.

El desarrollo de Cecilia 2b abre múltiples líneas de trabajo para el futuro. Es fundamental ampliar y diversificar el corpus de entrenamiento, integrando fuentes adicionales que reflejen una mayor variedad de registros, géneros y contextos de uso del español cubano. Asimismo, el entrenamiento con más épocas y el ajuste fino en tareas específicas permitirán mejorar la especialización y robustez del modelo, maximizando su utilidad en aplicaciones concretas.

Un objetivo clave será la creación de benchmarks y corpus de instrucciones diseñados específicamente para evaluar la comprensión y generación en español cubano, lo que permitirá medir de forma precisa el impacto y la ventaja competitiva de Cecilia en tareas relevantes para usuarios y desarrolladores locales. Finalmente, la integración de Cecilia en aplicaciones reales—como asistentes virtuales, sistemas de soporte educativo, o herramientas de procesamiento documental— será esencial para validar su utilidad práctica y fomentar su adopción en la sociedad cubana.

El desarrollo y despliegue de modelos de lenguaje regionales como Cecilia 2b plantea importantes consideraciones éticas. Es crucial asegurar que el uso de datos respete la privacidad, los derechos de autor y la diversidad cultural de las fuentes, evitando la reproducción o amplificación de sesgos presentes en los corpus de entrenamiento. La selección y documentación cuidadosa de los datos, así como la transparencia en los procesos de construcción y evaluación del modelo, son fundamentales para mitigar riesgos de homogeneización cultural y exclusión de voces minoritarias.

Además, el diseño participativo y la consulta con comunidades lingüísticas locales deben ser una prioridad, garantizando que las tecnologías desarrolladas respondan a las necesidades, valores y expectativas de sus usuarios, y contribuyan a una inteligencia artificial más justa, inclusiva y representativa.

4 Conclusiones

En este artículo se presenta Cecilia 2b, el primer modelo de lenguaje entrenado específicamente para el español cubano, construido mediante preentrenamiento continuo sobre la arquitectura Salamandra 2B sin modificar su estructura ni tokenizador. Cecilia utiliza un corpus diverso y representativo de la cultura y sociedad cubanas, compuesto por más de 296 mil documentos y cerca de mil millones de tokens, abarcando prensa, literatura, legislación y recursos enciclopédicos. El modelo se evaluó en una batería de tareas estándar multilingües, mostrando una reducción promedio de solo 2.4% respecto al modelo base, lo que indica que la especialización en datos cubanos no compromete su robustez general. Si bien la evaluación específica en tareas propias del español cubano es aún una asignatura pendiente, Cecilia 2b constituye un paso inicial fundamental hacia la creación de tecnologías lingüísticas adaptadas a la realidad

cubana y sienta las bases para futuros desarrollos y aplicaciones en procesamiento de lenguaje natural regionalizado.

Referencias

- [1] Samuel Cahyawijaya et al. “Crowdsource, crawl, or generate? creating sea-vl, a multicultural vision-language dataset for southeast asia”. In: *arXiv preprint arXiv:2503.07920* (2025).
- [2] Arnav Chavan et al. “One-for-all: Generalized lora for parameter-efficient fine-tuning”. In: *arXiv preprint arXiv:2306.07967* (2023).
- [3] Xin Dong et al. “Hymba: A hybrid-head architecture for small language models, 2024”. In: URL <https://arxiv.org/abs/2411.13676> ().
- [4] Lance Calvin Lim Gamboa and Mark Lee. “Filipino Benchmarks for Measuring Sexist and Homophobic Bias in Multilingual Language Models from Southeast Asia”. In: *arXiv preprint arXiv:2412.07303* (2024).
- [5] Aitor Gonzalez-Agirre et al. “Salamandra technical report”. In: *arXiv preprint arXiv:2502.08489* (2025).
- [6] Zhen Guo and Yining Hua. “Continuous training and fine-tuning for domain-specific language models in medical question answering”. In: *arXiv preprint arXiv:2311.00204* (2023).
- [7] Zijian Hu et al. “Fox-1 Technical Report”. In: *arXiv preprint arXiv:2411.05281* (2024).
- [8] Yoichi Ishibashi, Taro Yano, and Masafumi Oyamada. “Mining hidden thoughts from texts: Evaluating continual pretraining with synthetic data for llm reasoning”. In: *arXiv preprint arXiv:2505.10182* (2025).
- [9] Suramya Jadhav et al. “On Limitations of LLM as Annotator for Low Resource Languages”. In: *arXiv preprint arXiv:2411.17637* (2024).
- [10] Vladislav Lialin, Vijeta Deshpande, and Anna Rumshisky. “Scaling down to scale up: A guide to parameter-efficient fine-tuning”. In: *arXiv preprint arXiv:2303.15647* (2023).
- [11] Raymond Ng et al. “Sea-lion: Southeast asian languages in one network”. In: *arXiv preprint arXiv:2504.05747* (2025).
- [12] Kelechi Ogueji, Yuxin Zhu, and Jimmy Lin. “Small data? no problem! exploring the viability of pretrained multilingual language models for low-resourced languages”. In: *Proceedings of the 1st workshop on multilingual representation learning*. 2021, pp. 116–126.
- [13] Peng Shu et al. “Transcending language boundaries: Harnessing llms for low-resource language translation”. In: *arXiv preprint arXiv:2411.11295* (2024).
- [14] Atnafu Lambebo Tonja et al. “Inkubalm: A small language model for low-resource african languages”. In: *arXiv preprint arXiv:2408.17024* (2024).

- [15] Anh-Dung Vo et al. “Redwhale: An adapted korean llm through efficient continual pre-training”. In: *arXiv preprint arXiv:2408.11294* (2024).
- [16] Tongtong Wu et al. “Continual learning for large language models: A survey”. In: *arXiv preprint arXiv:2402.01364* (2024).
- [17] Yongyu Yan et al. “Af adapter: Continual pretraining for building chinese biomedical language model”. In: *2023 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. IEEE. 2023, pp. 953–957.
- [18] Yunzhi Yao et al. “Adapt-and-distill: Developing small, fast and effective pretrained language models for domains”. In: *arXiv preprint arXiv:2106.13474* (2021).