

## DESAFÍOS EN LA CONSTRUCCIÓN DE CORPUS PARALELOS EN ARAGONÉS: O CAMINO EN LA TRADUCCIÓN AUTOMÁTICA NEURONAL<sup>1</sup>

Alejandro PARDOS CALVO\*  
Universidad de Zaragoza

**RESUMEN** El proyecto *Traducción automática neuronal para las lenguas románicas de la península ibérica* (TAN-IBE), financiado por el Ministerio de Ciencia, Innovación y Universidades, tiene como objetivo entrenar sistemas de traducción automática basados en redes neuronales para siete lenguas románicas de la península ibérica, entre ellas el aragonés. Estos sistemas se entrenan mediante corpus paralelos, es decir, recopilaciones masivas de datos alineados que pueden utilizarse con distintos fines. En lenguas como el castellano, el portugués o el catalán —que cuentan con recursos claramente superiores—, la recopilación de esos materiales no plantea grandes dificultades. En cambio, para lenguas como el aragonés, el asturiano o el aranés, esta tarea resulta compleja debido a la escasez de recursos y de herramientas disponibles. En el caso del aragonés, a ello se suman otros obstáculos, como el bajo grado de digitalización, la considerable magnitud de los corpus dialectales o la diversidad de grafías y normas, entre otros. Este artículo analiza las dificultades específicas a las que se enfrentan los investigadores al recopilar o construir corpus paralelos en lengua aragonesa, evalúa los recursos existentes y propone posibles soluciones y estrategias para superar esos obstáculos.

**PALABRAS CLAVE** Lengua aragonesa. Corpus paralelos. Traducción automática neuronal. Entrenamiento de sistemas de traducción. *Transfer learning*. Sistemas multilingües.

**ABSTRACT** The *Neural machine translation for the Romance languages of the Iberian Peninsula* project (TAN-IBE), funded by the Ministry of Science, Innovation and Universities, aims to train neural network-based machine translation systems for seven Romance languages of the Iberian Peninsula, including Aragonese. These systems are trained using parallel corpora, that is, large collections of aligned data used for various purposes. For languages such as Spanish, Portuguese, or Catalan (which have indisputably stronger resources), data collection does not pose major difficulties. However, for languages like Aragonese, Asturian, or Aranese, it remains a demanding task due to the scarcity of available resources and tools. In the specific case of Aragonese, additional challenges arise, including

---

\* a.pardos@unizar.es

<sup>1</sup> Comunicación presentada en la IX Trobada d'Estudis e Rechiras Arredol d'a Luenga Aragonesa e a suya Literatura (Uesca, 3-5 d'otubre de 2024).

limited digitization, the considerable size of dialectal corpora, and the diversity of spellings and norms, among other factors. This paper analyzes the specific difficulties researchers face when collecting or constructing parallel corpora in Aragonese, assesses existing resources, and outlines potential solutions and strategies to overcome those obstacles.

**KEYWORDS** Aragonese language. Parallel corpora. Training of machine translation systems. Transfer learning. Multilingual systems.

**RÉSUMÉ** Le projet *Traduction automatique neuronale pour les langues romanes de la péninsule ibérique* (TAN-IBE), financé par le ministère espagnol des Sciences, de l'Innovation et des Universités, a pour objectif de former des systèmes de traduction automatique basés sur des réseaux neuronaux pour sept langues romanes de la péninsule ibérique, dont l'aragonais. Ces systèmes sont entraînés à l'aide de corpus parallèles, c'est-à-dire de vastes compilations de données alignées pouvant être utilisées à différentes fins. Dans des langues telles que l'espagnol, le portugais ou le catalan —qui disposent de ressources nettement supérieures—, la collecte de ces matériaux ne pose pas de grandes difficultés. En revanche, pour des langues comme l'aragonais, l'asturien ou l'aranais, cette tâche s'avère complexe en raison de la rareté des ressources et des outils disponibles. Dans le cas de l'aragonais, d'autres obstacles s'y ajoutent, tels que le faible degré de numérisation, l'ampleur considérable des corpus dialectaux ou la diversité des orthographes et des normes, entre autres. Cet article analyse les difficultés spécifiques auxquelles sont confrontés les chercheurs lors de la collecte ou de la construction de corpus parallèles en langue aragonaise, évalue les ressources existantes et propose des solutions et des stratégies possibles pour surmonter ces obstacles.

**MOTS CLÉS** Langue aragonaise. Corpus parallèles. Traduction automatique neuronale. Formation des systèmes de traduction. *Transfer learning*. Systèmes multilingues.

A tradución automática neuronal (TAN), alazetata en retes neuronals fundos, ha surtito en os zaguers años como una teunolochía abanzata ta traduzir entre difereus idiomas con buena precisión e fluidez. En contimparanza con otros sistemas de tradución, como aquers que s'alazetan en reglas u en frases, os modelos de TAN emplegan un aprendizache fundo ta analizar grans cantidaz de testos, asimilando patrons complexos en as estruturas gramaticals e semanticas d'as luengas. Ista enfiladura permite que os sistemas de tradución automática replequen matizes, senditos e contestos que se perdeban en otros metodos anteriores, chenerando asinas traduzioni más naturais e adecuatas ta ra situçaziòn (Genabith, 2020).

A más gran parti d'o desarrollo e entrenamiento de modelos de TAN s'ha zentrato en luengas mayoritarias como l'inglés, o castellano, o francés u o chino, en do sobrexen os recursos lingüísticos e os corpus paralelos de calidá e no bi ha problemas de financiación. Manimenos, ista enfiladura dixa de banda muitas atras luengas minoritarias que, á ormino, no poseen prous datos de calidá ta entrenar ditos modelos. Por tanto, os sistemas de tradución disponibles no consiguen replecar as características distintibas d'istas luengas, perpetuando ra suya imbisibilidad dichital e contrebuyindo á ra suya marguinaziòn en o contesto teunolochico (Oliver, 2020a). Ta ras luengas minorizadas como l'aragonés, a TAN s'amuestra alazetal ta ra suya preserbaziòn e ra suya promociòn.

L'aragonés, con una comunidá de fablans reducita e una presenzia encara limitata en os meyo de comunicaziòn, asinas como en os entornos dichitals, se troba

en una situación de vulnerabilidad sociolingüística que incrementa el riesgo de desplazamiento e, a largo plazo, de desaparición (Eito *et alii*, 2025). Frente a esto, el desarrollo de modelos de procesamiento del lenguaje natural entrenados de propósito para el aragonés abre una nueva línea de trabajo con potencial impacto en diferentes ámbitos: (i) puede mejorar la disponibilidad de herramientas de traducción y de producción textual, (ii) fortalecer la visibilidad y el empleo de la lengua en espacios digitales de comunicación cotidiana (como en redes sociales y blogs) y (iii) facilitar la generación de recursos didácticos y materiales de enseñanza y aprendizaje. Asimismo, la existencia de unos avances lingüístico-tecnológicos en aragonés puede contribuir a su incorporación en ámbitos institucionales y administrativos, con efectos positivos sobre la legitimidad percibida de la lengua y a su funcionalidad en registros formales (Baxter, 2017; Busquets, 2020).

Tal que la traducción automática neuronal pueda usufructuar todos sus beneficios a las lenguas minoritarias, es necesario poder contar con recursos lingüísticos adecuados, estando los corpus paralelos uno de los elementos principales en este proceso. Por *corpus paralelos* entendemos un conjunto de textos alineados en dos o más lenguas de cada segmento (trozo) en una lengua tiene una traducción equisivalente en la otra lengua (Gómez, 2005). Estos segmentos proporcionan datos comparables que ayudan a los modelos a replicar las relaciones lingüísticas entre las diferentes lenguas. Estos modelos aprenden, de forma organizada, cómo se traducen las palabras y las expresiones en determinados contextos y detectan patrones de correspondencia y equibalancia entre las lenguas en cuestión. Esencialmente, un corpus paralelo actúa como una base de datos bilingüe (o multilingüe) en que cada trozo o segmento tiene su propia traducción correspondiente en la otra lengua, lo que facilita la obtención de reglas de traducción de los modelos de TAN (Castillo, 2011). En lo que respecta a la temática, estos gozan de tener una gran variedad de géneros y estilos (literatura, noticias, documentación técnica, etc.), proporcionando así un abanico de registros y contextos de empleo, lo que constituye una representación de la variabilidad lingüística muy amplia.

Como se explicaba en el apartado anterior, el valor de la importancia de los corpus paralelos en el desarrollo de la traducción automática neuronal, no es el mismo para todas las lenguas, ya que algunas cuentan con una misma cantidad de recursos. Mientras que en las lenguas mayoritarias estos corpus son amplios y cubren gran variedad de contextos y géneros, en las lenguas minoritarias como el aragonés la realidad es bien diferente. La falta de corpus paralelos de calidad representa un obstáculo en la creación de modelos de TAN que sean precisos, dato que limita las posibilidades de traducir esta lengua fielmente y funcionalmente en diferentes contextos. En el presente artículo se analizará la situación de el aragonés en lo que respecta a la disponibilidad y la calidad de los corpus paralelos, ratiando de la experiencia recogida en el proyecto TAN-IBE (*Traducción automática neuronal para las lenguas románicas de la península ibérica*) (Oliver, 2020b).

Desde el inicio de mayo de 2022 hasta mayo de 2023, en el proyecto TAN-IBE, coordinado por la Universitat Oberta de Catalunya y financiado por el Ministerio de Ciencia, Innovación y Universidades, se trabajó en la recopilación, el análisis, la edición y la gestión

de corpus paralelos en aragonés (asinas como en as atras luengas romanicas minoritarias que fan parti d'ó proyeuto), prenendo ro castellano como luenga de referencia. Una güellada inicial dixaba alufrar que l'aragonés partiba d'una clara desbentalla fren á ra resta d'as luengas romanicas d'a peninsula iberica. Seguntes a web Opus Corpus, referencia clau en a replega de datos lingüísticos, l'aragonés se posiciona como a luenga romanica peninsular con menos segmentos disponibles, nomás por denzima d'ó mirandés, que en beras fa parti d'ó diasistema asturleonés, pero dispone d'un repositorio propio. Con alto u baxo sisanta mil segmentos disponibles, se troba muito luen de luengas minorizatas como l'oczitano, que cuenta con dos millons de segmentos, u l'asturiano, que, á penar d'a suya situgazió, tamién delicata, dispone de más de treze millons. Goyan de millor salú encara atras luengas, como ro catalán, o gallego u, profes, o portugués.

Portugués	482 900 000
Catalán	163 700 000
Gallego	67 900 000
Asturiano	13 300 000
Aranés (oczitano)	2 000 000
Aragonés	60 700
Mirandés	9 200

*Tabla 1. Lumero de segmentos d'os corpus paralelos entre o castellano e ra resta d'as luengas disponibles en a colezió d'Opus Corpus. (Fuen: Elaborazió propia)*

O prozeso de replega de corpus paralelos s'enzetó de traza organizata e con criterios esclateros ta guarenziar una obtenzió de datos efeutiba. Asinas, s'establioron tres linias de treballu estratechicas: se prebó de coordinar esfuerzos con as diferens instituzions poseedoras de material bilingüe; s'optó por a replega indibidual e por a creyazió de nuebos corpus paralelos por parti d'ó equipo rechirador; se prozedió á ra clasificazió d'os datos en dembas tematicas que respondesen á ros diferens contestos d'emplego.

En un primer inte se contactó con as instituzions e as entidaz prenzipals de l'ambito de l'aragonés, encluyendo-bi editorials, asoziacions culturals, prensa, meyos de comunicazió e creyadors de contenito, entre otros, con l'obchetibo d'identificar posibles fuens de corpus paralelos aragonés-castellano. As respuestas estioron bariatias: bellas entidaz colabororon de traza favorable, atras refusoron a suya partezipazió u determinoron de no responder. En iste causo, cal acobaltar o compromís que adotó ra Direzió Cheneral de Politica Lingüística, que, de camín, furnió

<sup>1</sup> Como no bi'n eba de corpus en aranés, emos establito ras zifras de l'oczitano.

una gran variedad de materiales, tanto bilingües como monolingües, que constituyeron la base temática más diversa de todo el proceso de replega. Además del contenido textual, la institución utilizó otros recursos de gran calidad, como el contenido de *l'Aragonario*, diccionario web castellano-aragonés, e incluso miles de lexemas que eran parados para puyar en Wikidata, enlazando así los recursos disponibles para desarrollar el TAN en esta lengua.

De vez, el personal reclutador procedió a la replega directa de corpus paralelos obtenidos de otras fuentes. Se adquirieron recursos como novelas, libros de diferentes temáticas o materiales didácticos digitalizados, páginas web bilingües, contenidos de blogs e redes sociales, artículos de prensa, revistas académicas, etcétera. Como novedad cabe mencionar la incorporación de todos los artículos de Wikipedia disponibles en aragonés, que constituyen buena parte de los recursos digitales, encara que en este caso presentan inconsistencias en cuanto a calidad o equibalancia, dato que muchos son estados generados automáticamente desde el castellano u desde otras lenguas. Para compensar estas limitaciones y reforzar las bases de datos se optó también por la traducción de nuevos corpus disponibles en castellano, como el *Flores-200*, que suministró miles de nuevos segmentos. A la revisión de estas aportaciones por parte de los expertos en cada una de las lenguas se ofrece una mayor precisión lingüística e mejora la calidad general de los datos disponibles.

Español	402 430
Portugués	429 730
Catalán	133 214
Gallego	39 627
Asturiano	11 734
Aranés (occitano)	14 584
Aragonés	10 552
Mirandés	2 000

Tabla 2. Número de artículos de Wikipedia disponibles en cada lengua.  
(Fuente: Elaboración propia)

Los datos replegados estuvieron organizados inicialmente en dos grandes categorías: corpus bilingües y corpus monolingües. Como se ha indicado antes, los corpus bilingües son ideales para el proyecto, dato que permiten el entrenamiento directo de los sistemas de traducción automática suministrando ejemplos esclarecidos de equibalancias frase a frase entre aragonés y castellano. Los corpus monolingües, encara que no contribuyen directamente a la creación de traducciones, constituyen un recurso de gran calidad para el análisis de patrones de gramática, léxico y estilo propios de aragonés, o que complementa la capacidad de o sistema para captar y generar contenidos más naturales en esta lengua (Sennrich, Haddow e Birch, 2015). Por eso se determinó de replegar corpus monolingües como complemento esencial para todas las lenguas del proyecto. Para

maximizar l'aplicabilidad d'o corpus replegado se lebó a cabo una clasificación temática que permitió guarenziar en buena medida la diversidad contextual. Así, se obtuvieron corpus en las áreas de derecho, administración, literatura, historia e cultura, educación, naturaleza e medioambiente.

Durante el proceso surgieron varios problemas. Muchas de ellas se podían atribuir al contexto general de las lenguas minoritarias; sin embargo, se observaron algunos problemas que afectan de forma específica al aragonés. Continúan presentándose los principales problemas identificados en esta fase del proyecto.

1. **Carencia de textos bilingües.** La propia inexistencia de textos bilingües aragonés-castellano representa el mayor problema. Del total de corpus paralelos, los corpus paralelos representan sólo el 47,62 %, mientras que los corpus monolingües representan el 52,38 %. Sin una base sólida de corpus paralelos con correspondencias claras, los modelos de TAN enfrentan graves problemas.
2. **Estandarización lingüística.** Una proporción significativa del total de los textos disponibles corresponde a ras diferenciales dialectales del aragonés. De la misma forma, la falta de un modelo lingüístico unificado, empleado de manera consensuada, dificulta la creación de corpus homogéneos. Esto condiciona tanto la alineación de textos entre las dos lenguas como la consistencia de las traducciones obtenidas.
3. **Diferencias combinatorias ortográficas.** La inconsistencia en la escritura del aragonés, resultante de la coexistencia de varias combinatorias ortográficas, dificulta nuevamente la creación de corpus uniformes. Esta variabilidad genera confusión en los sistemas de alineación y reduce la fiabilidad de los modelos de TAN, que requieren de una ortografía consistente para maximizar la precisión en el procesamiento de datos.
4. **Falta de digitalización.** Muchos de los textos escritos en aragonés aún no son digitalizados o se encuentran en formatos difícilmente procesables, mientras que la replega automatizada y el procesamiento de datos mediante software son aspectos esenciales para entrenar modelos de TAN de forma eficiente.
5. **Limitación de recursos web.** Aún cuando el aragonés siga incrementando su presencia en Internet durante los próximos años, esto continúa estando limitado. Esta limitación restringe el acceso a potenciales fuentes de gran calidad en el ámbito de contenidos actualizados e adaptados al contexto digital, representando, así mismo, una gran barrera de acceso.
6. **Calidad de las traducciones.** Las traducciones disponibles en aragonés varían mucho en cuanto a calidad. Se observa una gran diferencia, por ejemplo, entre las traducciones elaboradas durante las primeras décadas de la restauración de la lengua y las más actuales, ya que en esta última etapa las tendencias de empleo

encara no yeran definitas; antiparti, a luenga se trobaba en un inte d'adautazión á ras nuevas nezesidaz. Con tot e con ixo, bi ha muitas traduzions, fueras d'o rechistro emplegato, que no poderban estar consideratas material confitable, dato que son traduzions no professionals, reyalizadas por presonas afizionatas, e se trata d'un proyeuto que requiere prezisión.

7. **Limitazión en a variedá d'os rechistros e dembas tematicas.** L'aragonés ha desarrollato amplamén a suya literatura en as zagueras decadas. Manimenos, bi ha una gran manca en a produzién e a traduzión de testos teunicos, zientificos e legals, o que reduce l'aplicabilidá d'os modelos de TAN en contestos formals e espezializatos e dixta ra luenga en una posizién d'inferioridá en contimparanza con atras luengas más consolidatas.
8. **Sobrefaxo de recursos lexicograficos.** Existen lumerosos recursos lexicograficos, como dizionarios e replegas lexicas, que fan buena onra ta establir equibalenzias prezisas de parolas e ta enamplar o conoximiento d'o bocabulario disponible en dita luenga. Manimenos, istos recursos, fren á otros corpus d'uso prautico, restrinchen o entrenamiento de modelos en contestos reyals d'emplego d'a luenga.
9. **Restrizions d'azeso e dreitos d'autor.** Muitos materials son amparatos por dreitos d'autor, o que limita ra suya disponibilidá e ra suya incorporazién como corpus monolingües u paralelos. Á isto cal adibir a disparidá en a boluntá de colaborazién por parti d'os diferens autors implicatos, que incrementó ras dificultaz de replega.

Dimpués d'analizar con ficazio as dificultaz asoziatas á ra creyazién de corpus paralelos ta l'aragonés, cal planteyar azions que premitan sobrepuyar ista situ-gazién. Bi ha un dople obchetibo, que pasa por (i) enamplar a disponibilidá de corpus bilingües (e monolingües) de calidá e (ii) enfortir as ferramientas teunolochicas que premitan o desarrollo de modelos de TAN efziens.

1. **Replega de nuevos testos.** Como ye lochico, o medro d'o corpus aragonés pasarba en primer puesto por a identificazién e a replega de nuevos testos. Isto encluye a esplorazién de nuevas fuens literarias, dichitals e teunicas. A incorporazién de nuevos materials enampla ro rango de rechistros e estilos lingüísticos e permite entrenar modelos más zereños e adautables. Amás, a recopilazién sistematica contrebuye á ra creyazién d'una base de datos más representatiba d'os diferens emplegos de l'aragonés, que comprenda tanto ros contestos más informals como ros academicos u zientificos. Dato que muitas obras no son dichitalizadas, se suchiere priorizar a suya combersión enta formatos azesibles por meyo de teunicas de dichitalizazién como l'OCR (reconoximiento optico de caracteres).
2. **Colaborazions con autors e traductors locals.** O esito de cualesquier inziatiba de replega de corpus pende en gran mida d'a colaborazién autiba

con os prenzipals autors d'ó ecosistema lingüístico. En iste sendito, ta obtener materials ineditos u poco azesibles estarba alazetal establir alianzas con as editorials locals, asinas como con os autors e traductors más autibos. Iestas colaborazions enamplan o bolumen de datos disponibles e guarenzian un menimo de calidá d'os testos replegatos, ya que os professionals d'a luenga gosan estar millor paratos ta produzir traduzions coderens e fidels. Amás, a implicazió d'istos autors faborexe una mayor conzenzia e un mayor compromís en o desarrollo de ferramientas ta l'aragonés, creyando una sinerchia positiba entre ra lingüística e ra industria cultural.

3. **Plataformas de *crowdsourcing* ta ra replega de testos.** O *crowdsourcing* constituye una soluzión efeutiba e economica ta xamplar os corpus en o caso d'as luengas minoritarias (Capdevila, 2012). Á trabiés de diferens plataformas dichitals se puede embrecar á ra poblazió ta que furna traduzions, transcrizions u rebisions lingüísticas. S'espleita asinas a colaborazió comunitaria, pro útil ta luengas como l'aragonés, que tiene recursos raditos, pero cuenta con una comunidá de fablans autiba que puede contreyir en gran mida. Manimenos, cal implementar mecanismos de bali-dazió ta guarenzian una calidá e una coderenzia menimas en os materials furnitos. Una estratechia complementaria pasarba por establir combenios de colaborazió con instituzions d'educazió superior, como a Unibersidá de Zaragoza, que imparte un Diploma d'Espezializazió en Filolochía Aragonesa. Ista mena de colaborazions permite canaliar a esperenzia academica e guarenzia un fluxo contino de materials de calidá menima, elaboratos baxo supervisió dozén e con criterios filolochicos zereños, mantenendo a sostenibilidá d'os esfuerzos de replega á ro largo.
4. **Traduzió de testos teunicos, zientificos e legals.** A traduzió d'ista mena de materials ye alazetal ta enriqueir os corpus paralelos de l'aragonés, ya que istos rechistros espezializatos no gosan estar guaire representatos en o contesto d'as luengas minoritarias. Por exemplo, en l'ambito almenistratibo, a Comunidá Autonoma d'Aragón chenera un buen influmen de decumentazió legal e teunica como decretos, ordenanzas, resoluzions y autas que poderban traduzir-sen enta l'aragonés. Tamién, en l'ambito pribato, poderban redautar-se contratos estándar u reglamentos que reflexen os diferens conzeutos churidicos en amas luengas. Suzede ro mesmo con otros seutors, como l'agricultura, a chestió ambiental u ra incheniería, dembas relebans en a economía aragonesa, que producen manuals, guías e estudios teunicos que poderban estar traduzitos á l'aragonés. Finalmén, se poderban encluyir testos binclatos á ro patrimonio cultural e l'autibidá turistica. Materials como guías turisticas, catalogos, descrizió de rutas u folletos intormatibos pueden alportar terminolochía espezializata de bez que posizian l'aragonés como luenga de cultura e conoximiento en un contexto d'interpretazió patrimonial propia (Pardos e Sanagustín, 2022).

5. **Transcripción e traducción de registros orales.** L'aragonés tiene una tradición oral pro ampla que, bien emplegata, poderba constituyir un recurso de gran balura ta ros corpus lingüísticos. A transcripción e traducción de registros orals como entrebistas, relatos, cançons u otras composicions furnirba un buen exemplo d'emplego d'a luenga en contestos cutianos.
6. **Desarrollo de correutors ortograficos, lematizadors e analizadors espezi-ficos ta l'aragonés.** A creyazió de ferramientas lingüísticas espezi-ficas ye alazetal ta prozesar e normalizar os textos replegatos. Isto permite estandarizar a ortografía, resolber as inconsistenzias identificatas e estruturar os textos en un formato adecuado ta ra suya incorporazió como corpus paralelos. Estarba espezialmén útil, ya que l'aragonés cuenta con diferens combenzions ortograficas e barians dialeutals que, en a practica, pueden comprometer a unificazió d'os datos lingüísticos.
7. **Aplicazió de teunicas d'aprendizache que profiten datos de luengas romanicas zercanas.** L'aragonés, como luenga romanica, poseye carauteristicas estruturals e lesicas comuns con otras luengas d'a suya familia, como ro catalán. Ista prosimidá lingüística puede estar espleitata meyán teunicas de transferencia d'aprendizache (*transfer learning*) e modelato multilingüe (Nguyen e Chiang, 2017). D'ista traza, se pueden entrenar modelos multilingües que encluyan grans corpus de catalán u gallego ta millorar o rendimiento d'a traducción automatica en aragonés.

O desarrollo de corpus paralelos ta l'aragonés ye una fayena estrategica de gran complexidá, que traszende l'ambito tecnico e se fica en o nucleo d'os esfuerzos por rebitalizar ista luenga minoritaria. En o presén articlo s'han identificato tanto ras limitacions estruturals que enfrontina iste proceso como ras oportunidaz de millora que poderban abordar ditas carenzias. Entre ros prenzipals barraches se puede acobaltar a escasez de textos bilingües e ra manca d'estandarizazió lingüística, que dificultan l'aliniazió prezisa de materials e comprometen a consistenzia d'as traduccions automaticas. Istos desafíos son, en muitos causos, inerens á ra reyalidá d'as luengas minoritarias, pero se beyen agrabatos por custions que afeutan particularmén á l'aragonés, como ra coesistenza de diferens combenzions ortograficas u ra presenzia radita en o entorno dichital. Á penar de tot ixo, os abanzas reziens en ferramientas dichitales e metodolochías de prozesamiento de luengache ubren nuevos caminos enta ra recuperazió e ra superazió d'istas dificultaz.

As estrategias proposatas miran d'empentar nuevos paradigmas en as trazas de replegar, prozesar e emplegar datos lingüísticos en aragonés. As inziatibas más importans pasarban por a colaborazió con instituzions culturals e academicas, a creyazió de plataformas de *crowdsourcing* e l'aplicazió de teunicas d'aprendizache transferito dende otras luengas romanicas zercanas. Amás, a dibersificazió tematica e ra incorporazió de registros espezializatos como ros textos zientificos, legals u educatibos resultan alazetals ta guarenziar que ros modelos puedan responder á

ras nezesidaz practicas e cutianas d'os fablans e ras demandas de ziertos seutors, de bez que premiten enfortir o prestichio d'a luenga en contestos contemporanios.

L'aragonés difizilmén podrá abanzar en rebitalizazi3n e normalizazi3n sin corpus paralelos. En un entorno cada begata m3s dichitalizado, constituyen una ferramienta que premite treballar con rigor e continuid3. O esito pender3, en gran mida, d'a coordinazi3n entre autors locals, instituzions e espertos en prozesamiento d'o luengache natural, como se ye fendo dende ro proyeuto TAN-IBE. Empentar ista mena de propuestas premitir3 consolidar l'aragonés como una luenga de raso funzional en o entorno dichital e ficar3 ros alavez ta futuras autuazi3ns destinadas 3 ro suyo desarrollo en o sieglo XXI.

## BIBLIOGRAFÍA

- Baxter, Robert Neal (2017), «The importance of interpreter training for minority languages: an analytical overview of the co-official languages in Spain», *Quaderns*, 24, pp. 151-177.
- Busquets, Blanca (2020), «Per què és important la traducció de llengües minoritàries?», en *Arts i Humanitats*, blog de l'Universitat Oberta de Catalunya <<https://n9.cl/kq2fh>> [consulta: 15/5/2020].
- Capdevila Fernández, Cristian (2012), «Crowdsourcing i traducció / localització: una amenaça o una oportunitat?», *Revista Tradumàtica*, 10, pp. 237-243.
- Castillo Rodríguez, Cristina (2011), «La alineación de un corpus paralelo multilingüe: propuesta de fases para la didáctica de traducción especializada inversa», *Cadernos de Tradução*, 1 (27), pp. 117-142 <<https://dialnet.unirioja.es/descarga/articulo/4925266.pdf>>.
- Eito Mateo, Antonio, José Ángel Iranzo Sanz, Chaime Marcuello Servós e Alejandro Pardos Calvo (2025), *Charrando aragonés: la lengua aragonesa en su zona de uso predominante*, Zaragoza, PUZ.
- Genabith, Josef van (2020), «Neural Machine Translation», en Jörg Porsiel (ed.), *Maschinelle Übersetzung für Übersetzungsprofis: Sammelband*, Berlín, BDÜ.
- Gómez Guinovart, Xavier (2005), «Procesamiento y aplicaciones de los corpus paralelos», *Novática*, 175, pp. 50-54.
- Nguyen, Toan Q., e David Chiang (2017), «Transfer learning across low-resource, related languages for neural machine translation», en Greg Kondrak e Taro Watanabe (eds.), *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, Taipéi, Asian Federation of Natural Language Processing, pp. 296-301.
- Oliver, Antoni (2020a), «Human translation and machine translation: specificities, uses, advantages and disadvantages», *Linguapax Review*, 8, pp. 111-131 <<https://n9.cl/d3odn>>.
- (2020b), «Traducción automática para las lenguas románicas de la península ibérica», *Studia Romanica et Anglica Zagrabiensia*, 65, pp. 367-375.
- Pardos Calvo, Alejandro, y Victoria Sanagustín Fons (2022), «Turismo cultural y lengua aragonesa: un estudio de caso en el municipio de San Juan de Plan (valle de Chistau, Huesca, Aragón)», *ROTUR*, 16 (2), pp. 186-206 <[https://revistas.udc.es/index.php/rotur/article/view/rotur.2022.16.2.8995/g8995\\_pdf](https://revistas.udc.es/index.php/rotur/article/view/rotur.2022.16.2.8995/g8995_pdf)>.
- Sennrich, Rico, Barry Haddow e Alexandra Birch (2015), «Improving neural machine translation models with monolingual data», en Katrin Erk e Noah A. Smith (eds.), *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Berlín, Association for Computational Linguistics, pp. 86-96.