



L'horitzó d'Aina: estatus del projecte

Actualització sobre el desenvolupament dels models lingüístics, recursos i objectius d'AINA

19 de desembre del 2023



Recordem

OBJECTIUS

1

Proveir el català de la **infraestructura** necessària **per al desenvolupament d'aplicacions basades en IA/TL**, (assistents de veu, traductors automàtics, agents conversacionals, etc)

2

Fer que **la inclusió del català** a les aplicacions de IA/TL sigui **rendible i atractiva per a les empreses del sector**, tant a nivell local com global.

3

Aconseguir que el ciutadà de Catalunya pugui **participar en català** en el món digital **al mateix nivell que un parlant d'una llengua global**, com ara l'anglès o el castellà.



DEFINICIÓ i FONAMENTS

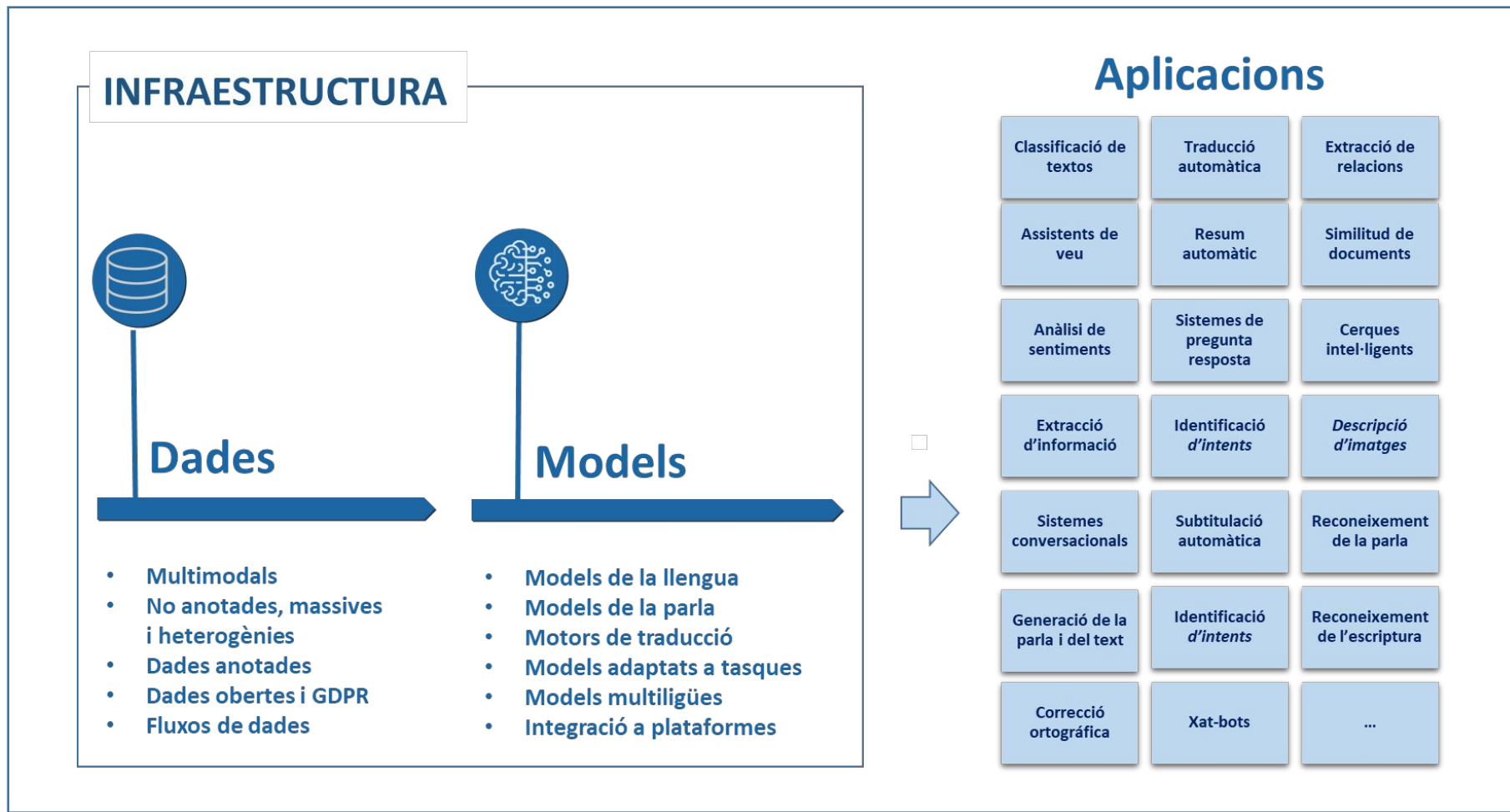
- AINA és essencialment **infraestructura** lingüística.
- El valor de les **dades**.
 - La tecnologia avança molt ràpidament però **les dades són persistents**.
 - Únicament des de la **iniciativa pública**, el català es pot garantir el subministrament de dades suficients.
 - Disposar de **dades de qualitat suficients és un actiu segur i de futur** que garanteix l'actualització de la tecnologia.

ESTRATÈGIA

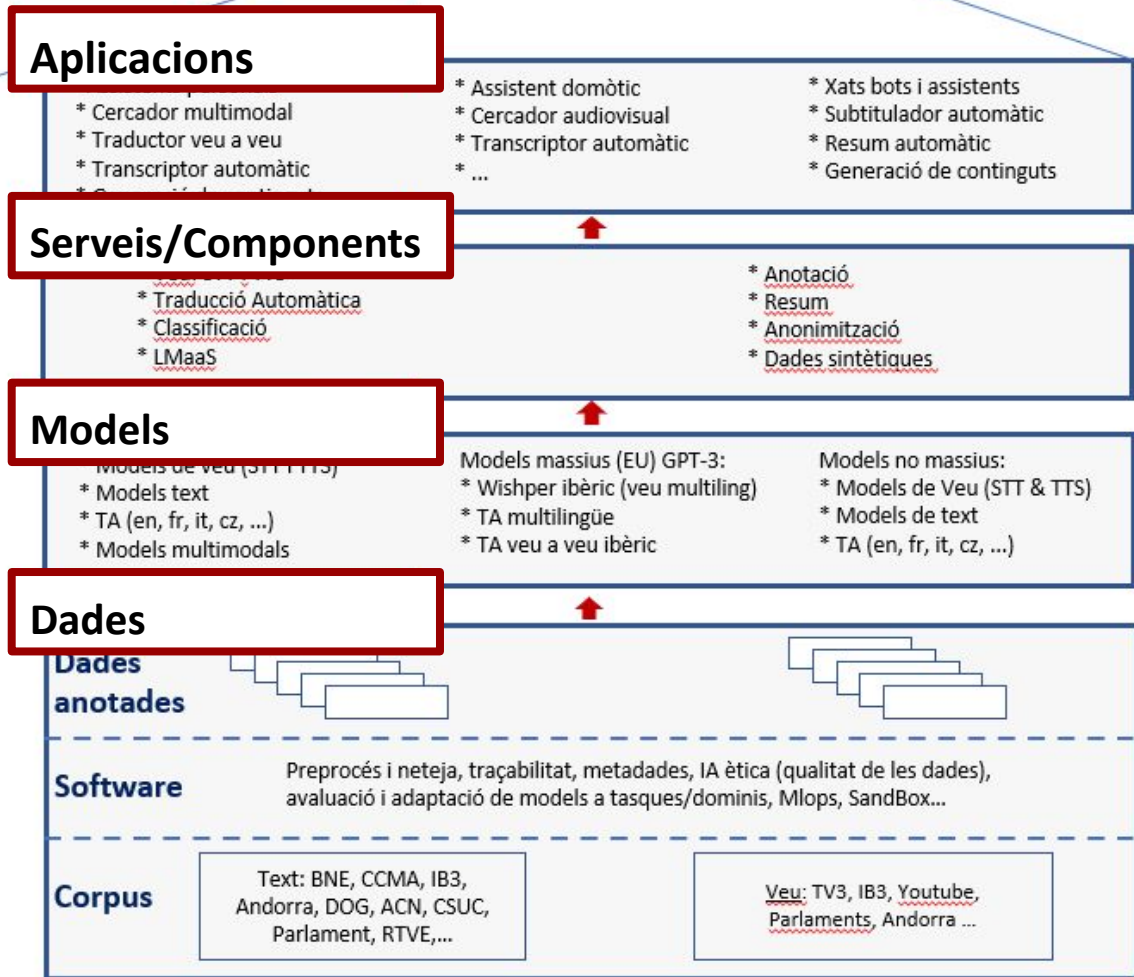
- AINA implementarà una infraestructura de recollida i neteja de dades amb la **implicació de grans actors**.
- AINA reaccionarà ràpidament als **avenços tecnològics** mitjançant la vigilància tecnològica.
- AINA detectarà i donarà resposta a **noves necessitats** de les empreses i de la societat mitjançant la vigilància sectorial i de mercat.

Barcelona, 15 de febrer del 2022

Recordem



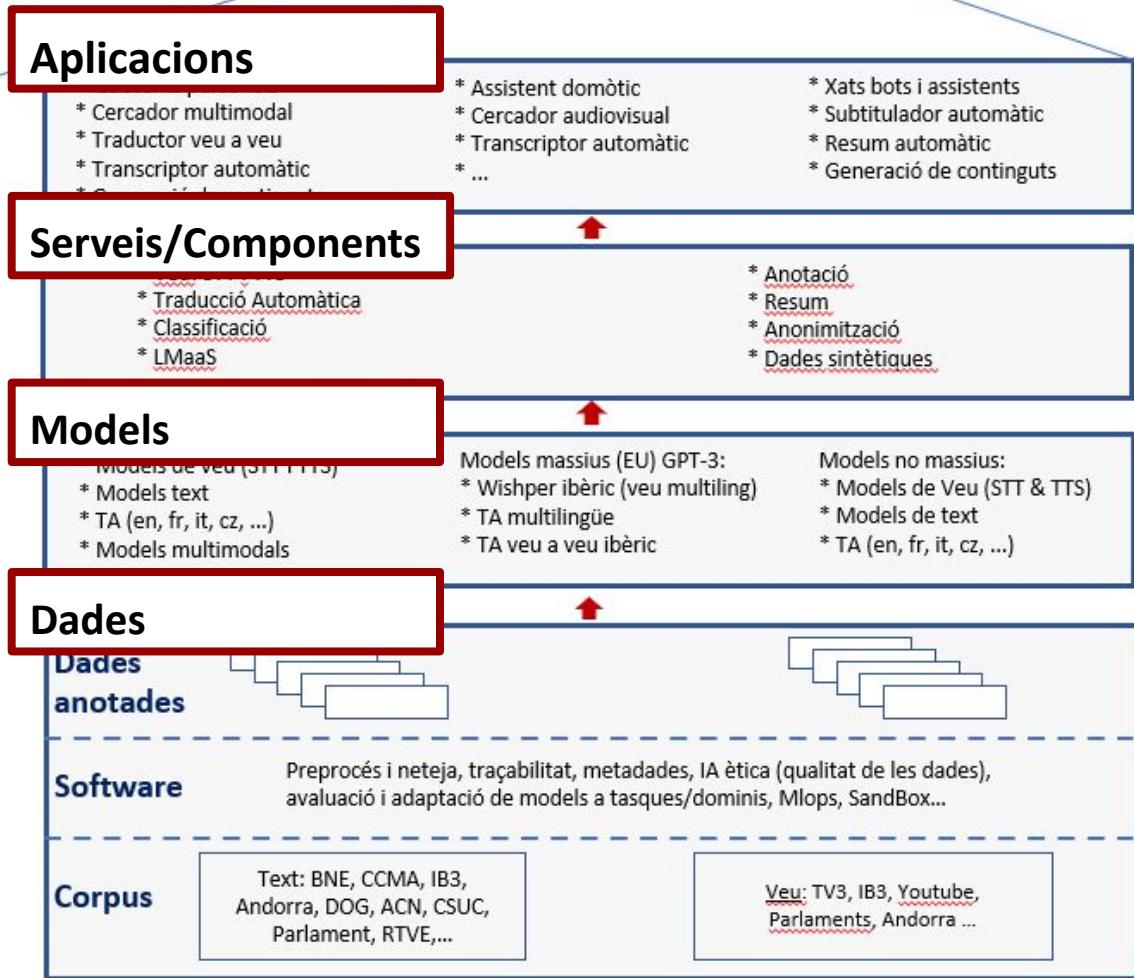
INFRAESTRUCTURA LINGÜÍSTICA IA- AINA



Provisió i processament de les dades lingüístiques

- Corpus textual
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades

INFRAESTRUCTURA LINGÜÍSTICA IA- AINA



Models

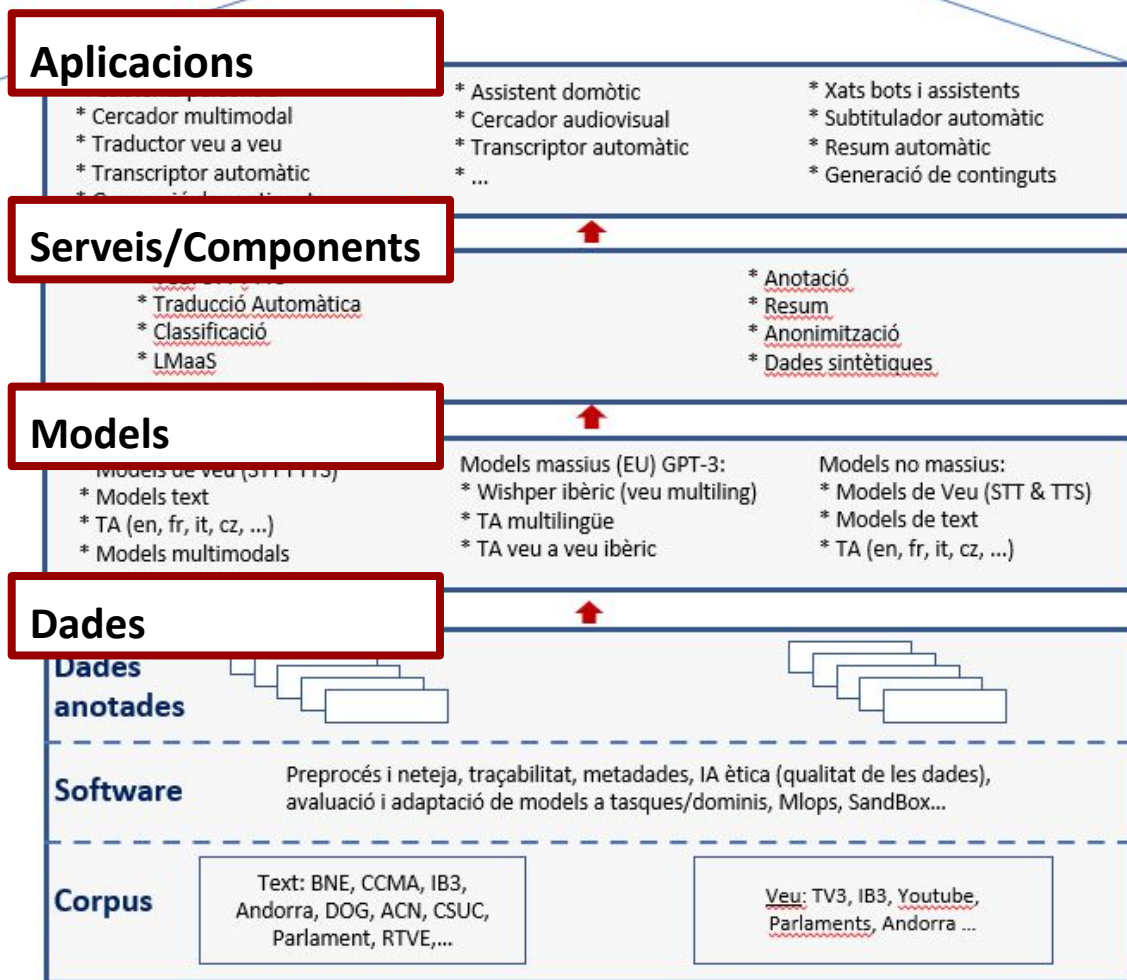
- Generació de **models** de llengua
- Tecnologies de la **parla**
- Tecnologies de la **traducció**

Eines d'avaluació i benchmarking La IA ètica

Provisió i processament de les dades lingüístiques

- Corpus textual
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades

INFRAESTRUCTURA LINGÜÍSTICA IA- AINA



Serveis lingüístics

- Accés i desplegaments
- Integració en frameworks i plataformes de referència.
- APIS, demosDemos

Models

- Generació de **models** de llengua
- Tecnologies de la **parla**
- Tecnologies de la **traducció**

Eines d'avaluació i benchmarking La IA ètica

Provisió i processament de les dades lingüístiques

- Corpus textual
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades



Provisió de dades

- **Corpus textual**
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades

V3. Corpus textual, pre-processat i publicat *

Corpus	v1 (GB)	v2 (GB)	v3 (GB)	Estat	Disponible a	Llicència
DOGC	0.78	0.78	0.92	v2	OPUS	CC0 4.0
Catalan Open Subtitles	0.02	0.02	0.02	v1	OPUS	Oberta
Catalan Oscar	4.00	4.00	27.10	v2	OSCAR	CC0 4.0
CaWaC	3.60	3.60	8.70	v2	CaWac	CC-BY-SA-3.0
Cat. General Crawling	2.50	2.50	5.70	v2	Zenodo	CC-BY 4.0
MaCoCu	-	-	11.00	Nou	MaCoCu	CC-BY 4.0
Viquipèdia	0.98	1.10	1.70	v3	-	CC-BY 4.0
ACN	0.45	0.45	0.45	v1	Zenodo	CC-BY-NC-ND 4.0
NacióDigital	-	0.45	0.45	v1		
VilaWeb	-	0.06	0.30	v2		
Catalan mc4	-	-	41.00	Nou	HuggingFace	CC-BY 4.0
Grup El Món	-	-	0.55	Nou	-	CC-BY 4.0
RacoCatalà	-	8.10	11.00	Nou	HuggingFace	CC-BY-NC 4.0
Altres	0.24	13.89	18.26	v3	-	CC-BY 4.0
Total	12.57	34.95	127.5			



* Totes aquestes dades es publicaran properament a [HuggingFace](#), sota el nom de CATalog 1.0, i constituïran el dataset de preentrenament de LLMs en català més gran fins ara.



Provisió de dades

- **Corpus textual**
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades

- **CATALOG 1.0:** és el dataset de preentrenament de LLMs més gran en català alliberat fins ara. A més, conté una gran varietat de fonts, amb un percentatge important de textos curats manualment, cosa que el diferencia enormement dels altres grans datasets publicats, que estan constituïts únicament per dades d'origen web.

Està format per **17.450.496.729 de paraules** (al voltant de **23 mil milions de tokens**) distribuïts en 34.314.510 documents.

- Accés: <https://huggingface.co/datasets/projecte-aina/CATalog>

The following pre-existing datasets

Media Groups	Academic & Book Repositories
<ul style="list-style-type: none">● OSCAR-2301● OSCAR-2201● CaText● MaCoCu-ca 1.0● caWaC● Colossal OSCAR 1.0● mC4	<ul style="list-style-type: none">● Tesis Doctorals en Xarxa (TDX)● Wikipedia● Project Gutenberg● Government Institutions● Parlament de Catalunya● Les Corts Valencianes● Diari Oficial de la Generalitat Valenciana● Butlletí Oficial de la Universitat d'Alacant



Provisió de dades

Creació de col·leccions de dades anotades per *fine tuning* i/o avaluació

- Corpus textual
- **Generació de dades anotades**
 - Generació d'instruccions
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades



CEIL: Corpus de 60.000 textos curts per a la identificació, classificació i vinculació d'entitats. Conté 9 tipus i 52 subtipus.

- Tasca: Identificació, classificació i vinculació d'entitats.
- Accés: <https://huggingface.co/datasets/projecte-aina/ceil>



CoQCat: Corpus de 6.000 paràgrafs anotats amb una conversa d'uns 15 torns de pregunta-resposta.

- Tasca: Pregunta resposta conversacional.
- Accés: <https://huggingface.co/datasets/projecte-aina/CoQCat>



CaSSA: Corpus de 6.400 ressenyes i missatges de fòrum anotats amb expressions de polaritat.

- Tasca: Anàlisi de sentiments.
- Accés: <https://huggingface.co/datasets/projecte-aina/CaSSA-catalan-structured-sentiment-analysis>



CaSET: Corpus de tuits anotats amb emocions i opinió. Conté 11.000 frases úniques sobre cinc temes controvertits, agrupades en 6.000 parells de frases.

- Tasca: Detecció d'emocions i opinió.
- Accés: <https://huggingface.co/datasets/projecte-aina/CaSET-catalan-stance-emotions-twitter>



CaSERa: Corpus de missatges de fòrum anotats amb emocions i opinió. Conté 15.782 frases úniques agrupades en 10.745 parells de frases.

- Tasca: Detecció d'emocions i opinió.
- Accés: <https://huggingface.co/datasets/projecte-aina/CaSERa-catalan-stance-emotions-raco>



Provisió de dades

- Corpus textual
- **Generació de dades anotades**
 - Generació d'instruccions
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades



CaBreu: Corpus de resums abstractius, extractius i extrems de 3.000 documents.

- Tasca: Resum de textos.
- Accés: <https://huggingface.co/datasets/projecte-aina/caBreu>



NLUCat: Corpus de 12.000 frases anotades segons el seu *intent* i els *slots* més rellevants, com a dades de suport per al desenvolupament d'assistents electrònics.

- Tasca: Detecció d'intencions.
- Accés: <https://huggingface.co/datasets/projecte-aina/NLUCat>



InToxiCAT: Corpus de 29.809 frases obtingudes de missatges de fòrums, anotades segons si són o no abusives.

- Tasca: Detecció de llenguatge abusiu.
- Accés: <https://huggingface.co/datasets/projecte-aina/InToxiCat>



COPA-ca: Traducció al català del corpus COPA, corpus de referència per al raonament causal. Inclou 1.000 instàncies, cadascuna de les quals es compon d'una premissa i dues hipòtesis (o alternatives).

- Tasca: Raonament causal.
- Accés: <https://huggingface.co/datasets/projecte-aina/COPA-ca>



PAWS-ca: Traducció al català del corpus PAWS per a la indentificació de paràfrasis. Conté 4.000 parells d'exemples traduïts per humans i 49.400 parells traduïts automàticament.

- Tasca: Paràfrasi.
- Accés: <https://huggingface.co/datasets/projecte-aina/PAWS-ca>



Provisió de dades

- Corpus textual
- **Generació de dades anotades**
 - Generació d'instruccions
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades



NOU

- **XNLI-ca**: Traducció al català del corpus XNLI per a l'avaluació de sistemes de comprensió de llenguatge cross-lingüe en tasques com ara la inferència del llenguatge natural. Conté 7.500 parells de frases.
 - Tasca: Implicació textual.
 - Accés: <https://huggingface.co/datasets/projecte-aina/xnli-ca>
- **WNLI-ca**: Traducció al català del corpus WNLI, que inclou 855 parells d'oracions, en què la primera frase conté una ambigüitat i la segona una possible interpretació d'aquesta.
 - Tasca: Implicació textual.
 - Accés: <https://huggingface.co/datasets/projecte-aina/wnli-ca>
- **Teca**: Corpus d'implicació textual. Conté 21.163 parells de premisses i hipòtesis, anotades segons la relació d'inferència que tenen (implicació, contradicció o neutre).
 - Tasca: Implicació textual
 - Accés: <https://huggingface.co/datasets/projecte-aina/teca>
- **TeCla v2**: Corpus de notícies en català per a tasques de classificació de textos multiclasse.
 - Tasca: Classificació de textos.
 - Accés: <https://huggingface.co/datasets/projecte-aina/tecla>
- **WikiCAT_ca**: Corpus català per a tasques de classificació temàtica de textos no periodístics.
 - Tasca: Classificació de documents.
 - Accés: https://huggingface.co/datasets/projecte-aina/WikiCAT_ca



Provisió de dades

- Corpus textual
- **Generació de dades anotades**
 - Generació d'instruccions
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades

- **CaWikiTC**: Corpus creat de manera automàtica a partir dels resums d'articles de la Viquipèdia i la categoria temàtica associada. Conté 21.002 textos classificats en 67 categories.
 - Tasca: Classificació de textos.
 - Accés: <https://huggingface.co/datasets/projecte-aina/CaWikiTC>
- **XQUAD-ca**: Traducció al català del dataset XQUAD, que consta d'un subconjunt de 240 paràgrafs i 1.190 parells de pregunta-resposta del dataset SQuAD v1.1
 - Tasca: Pregunta resposta multilingüe.
 - Accés: <https://huggingface.co/datasets/projecte-aina/xquad-ca>
- **VilaQuAD**: Corpus de parells preguntes/resposta sobre notícies. Conté 2.095 articles de notícies en català i d'1 a 5 preguntes amb la seva resposta per a cada fragment (o context).
 - Tasca: Pregunta resposta.
 - Accés: <https://huggingface.co/datasets/projecte-aina/vilaquad>
- **ViquiQuAD**: Corpus de parells preguntes/resposta sobre la Viquipèdia. Conté 3.111 contextos extrets d'un conjunt de 597 articles originals i d'1 a 5 preguntes amb la seva resposta per a cada fragment.
 - Tasca: Pregunta resposta.
 - Accés: <https://huggingface.co/datasets/projecte-aina/viquiquad>
- **CatalanQA**: Corpus de 21.426 parells preguntes/resposta sobre la Viquipèdia i notícies.
 - Tasca: Pregunta resposta.
 - Accés: <https://huggingface.co/datasets/projecte-aina/catalanqa>



Provisió de dades

- Corpus textual
- **Generació de dades anotades**
 - Generació d'instruccions
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades

- **CAT ManyNames:** Versió catalana del corpus *ManyNames*, orientat a models de Llenguatge i Visió. Consisteix en més de 23.000 imatges i les anotacions corresponents, traduïdes automàticament de l'anglès, més un conjunt de 1.072 imatges anotades a mà directament en català.
 - Tasca: Identificació d'imatges.
 - Accés: https://huggingface.co/datasets/projecte-aina/cat_manynames
- **XitXat:** Corpus de 950 converses de 10 dominis diferents entre xatbots i usuaris.
 - Tasca: Classificació d'*intents*, detecció d'entitats associades i entrenament/avaluació de sistemes conversacionals.
 - Accés: <https://zenodo.org/record/7276036#.Y2zMn4LMITU>
- **Parafraseja:** Corpus de 21.984 parells de frases anotades segons si són paràfrasis l'una de l'altra, o no.
 - Tasca: Paràfrasi.
 - Accés: https://huggingface.co/datasets/projecte-aina/Para_fraseja
- **STS-ca:** Corpus per a l'avaluació de la similitud textual semàntica.
 - Tasca: Similitud textual semàntica.
 - Accés: <https://huggingface.co/datasets/projecte-aina/sts-ca>
- **NoNiRes:** Anotació de les expressions de negació de 20.541 frases en català.
 - Tasca: Negació.
 - Accés: https://zenodo.org/record/7319487#.Y3S_uL7MLOs



Provisió de dades

- Corpus textual
- **Generació de dades anotades**
 - Generació d'instruccions
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades

- **GuiaCat:** Corpus de 5.750 ressenyes de restaurants en català de la plataforma GuiaCat. Cada ressenya té associada una valoració per servei, menjar, qualitat-preu i ambient, i una nota mitjana.
 - Tasca: Anàlisi de sentiments.
 - Accés: <https://huggingface.co/datasets/projecte-aina/GuiaCat>
- **ANCORA_ca v2:** Corpus d'entrenament de cadenes de processament, afegint la columna NER a la versió CONLLU de UD versió 9, per fer *multitask learning* dins de Spacy.
 - Tasca: Anotació d'entitats amb nom i dependències.
 - Accés: <https://doi.org/10.5281/zenodo.5036650>
- **AnCora-Ca-NER:** Corpus d'entrenament per a l'anotació d'entitats, basat en el corpus ANCORA de la UB.
 - Tasca: Anotació d'entitats.
 - Accés: <https://huggingface.co/datasets/projecte-aina/ancora-ca-ner>



Provisió de dades

Creació de col·leccions de dades d'instrucció de models

- Corpus textual
- Generació de dades anotades
 - **Generació d'instruccions**
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades

InstruCAT: Conversió de dades de datasets d'AINA per instruir LLMs en català:

216 mil instruccions per entrenament de LLMs en tasques “downstream”, com pregunta/resposta, resum, , parafrasi, classificació, etc.

Generades a partir dels següents **datasets:** caBreu, CatalanQA, CoQCat, GuiaCat, IntoxiCat, Parafraseja, PAWS-ca, sts-ca, WikiCat, CEIL, TECA, NLUcat.

Category	Number of instructions	%
ner	59410	25.24%
paraphrasis	34695	14.74%
text_classification	33393	14.19%
toxicity	29809	12.66%
qa	27427	11.65%
emotion_detection	18492	7.85%
phrase_generation	11873	5.04%
entailment_generation	6354	2.70%
sentiment_analysis	5750	2.44%
abstractive_summarization	2999	1.27%
extreme_summarization	2999	1.27%
entailment	2117	0.89%



<https://huggingface.co/datasets/BSC-LT/InstruCat>



Provisió de dades

- Corpus textual
- Generació de dades anotades
 - **Generació d'instruccions**
- Corpus de veu
 - Campanya de recollida de veu (CM)
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades



Generació d'un **dataset d'instruccions es-ca**: Dataset d'instruccions en castellà [[MentorES](#)] (i posterior traducció al català) per avaluar sistemes d'intel·ligència artificial. Inclou 10.175 tasques, distribuïdes en les següents categories: closed_qa, classification, open_qa, summarization, general_qa, information_extraction, brainstorming i creative_writing.



Traducció al català de **datasets d'instruccions**:

- Dolly
- OpenAssistant



Traducció al català de **datasets d'avaluació**:

- ARCChallenge i ARCEasy (raonament i coneixement bàsic sobre ciències).
- PIQA, (raonament lògic sobre el món físic).
- OpenBookQA (coneixements bàsics del món i certa capacitat de raonament).
- MGSM (raonament matemàtic).
- XStory Cloze (comprensió del llenguatge i raonament espacial i temporal propi de les narracions).



Provisió de dades



- Corpus textual
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - **Campanya recollida de veu CM**
 - Corpus de veu amb transcripció
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades

El català, entre les llengües 'top' en l'entrenament del nou model de traducció automàtica de Google

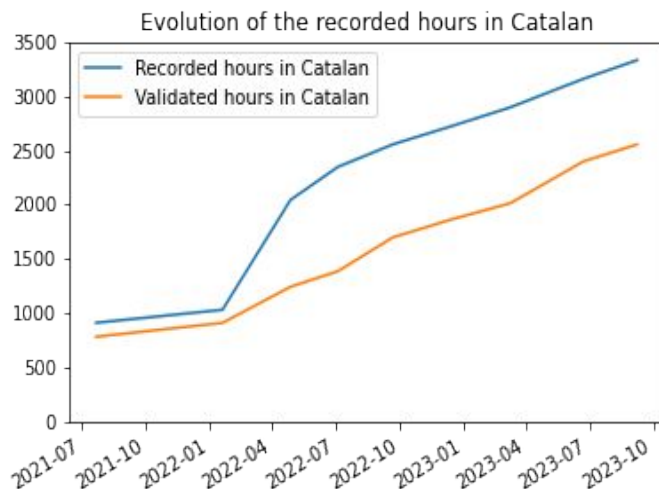
L'AudioPaLM, capaç de generar text i veu en multitud d'idiomes, es nodreix d'un repositori de veus amb aportacions de tot el món



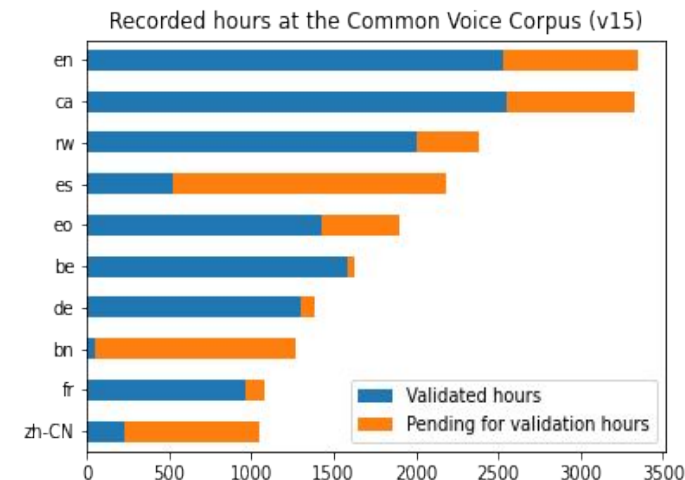
La recollida de dades mitjançant el **Common Voice** ha continuat estable:

- Des del gener de 2023 s'han recollit **més de 800 hores**.
- Des de l'inici del projecte s'han recollit **més de 2.500 hores**.
- Hi han col·laborat més de **35.000 voluntaris**

El català s'ha situat com a **primera llengua del corpus** pel que fa a **hores validades** i com a **segona llengua amb més hores enregistrades**, a molt poca distància de la primera



Evulució d'hores enregistrades i validades



Llengües amb més hores enregistrades a la v. 15

Common Voice
moz://a

<https://commonvoice.mozilla.org/ca/datasets>



Provisió de dades

- Corpus textual
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - Campanya recollida de veu CM
 - **Corpus de veu amb transcripció**
 - Corpus d'altres tasques de veu
- Corpus paral·lels per TA
- Estratègia de provisió de dades

- **ParlamentParla**: Corpus d'àudio de 611h d'enregistraments del Parlament. La preparació de la versió 3 està en marxa amb un mínim de 300h més d'enregistraments..
 - Tasca: Reconeixement de Parla.
 - Accés: https://huggingface.co/datasets/projecte-aina/parlament_parla



- **CCMA TV3**: Corpus d'àudio de 735h d'enregistraments de diversos programes de TV3
 - Tasca: Reconeixement de Parla
 - Accés: Per l'ús intern, s'està parlant amb la CCMA per alliberar-ho.



- **TTS cat**: Versions netes dels corpus de FESTCAT i OpenSLR69, específicament preparades per entrenar models neuronals de TTS
 - Tasca: Síntesi de la parla.
 - Accés: https://huggingface.co/datasets/projecte-aina/festcat_trimmed_denoised/
 - Accés: <https://huggingface.co/datasets/projecte-aina/openslr-slr69-ca-trimmed-denoised/>

ILENIA
IMPULSO DE LAS LENGUAS
EN LA INTELIGENCIA ARTIFICIAL



- **Corts Valencianes (ILENIA)**: Corpus d'àudio amb més de 300h d'enregistraments.

- Tasca: Reconeixement de la Parla.



- **IB3 (ILENIA)**: Corpus d'àudio amb 30h de programació d'IB3.

- Tasca: Reconeixement de la Parla.



Provisió de dades

- Corpus textual
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - Campanya recollida de veu CM
 - Corpus de veu amb transcripció
 - **Corpus d'altres tasques de veu**
- Corpus paral·lels per TA
- Estratègia de provisió de dades



NOU *Transcripció Fonètica*: Transcripció fonètica multi dialectal pel català. Inclou 4 dialectes, 160 frases cada un: central, nord-occidental, valencià, balear.

- Tasca: Transcripció fonètica.
- Accés: <https://huggingface.co/datasets/projecte-aina/4catat>



NOU *Annotated Catalan CV*: Anotació del corpus Common Voice amb gènere i dialecte dels parlants. En col·laboració amb el CLiC (Centre de Llenguatge i Computació) - UB

- Tasca: Identificació de gènere, identificació de dialecte
- Accés: https://huggingface.co/datasets/projecte-aina/annotated_catalan_common_voice

[ice](#)

sentence string · lengths	transcription string · lengths
 26 131	 31 146
Ses coses importants són ses que no ho semblen	səs k'ɔzəz import'an s'on səs kə n'o w s'əmblən
Vaig obrir sa caixa que temps enrere vaig omplir de records des meu germà	v'adʒ obr'i sə k'afə kə t'enz ənɾ'erə v'adʒ ompl'i ðə rək'ɔr ðəz m'ew zərm'a
Passejàvem ara pes estrets carrers que tants de cops havíem recorregut feia anys	pəsəðʒ'avəm 'arə pəðz əstr'ək kər'es kə t'aŋ də k'ɔðz əv'ian rəkərəv'uf f'əjə 'ajns
Te deix amor la mar com a penyora	tə ð'eʒ əm'o lə m'a k'om ə pen'orə
Setze exercicis d'un jutge	s'edʒ əðʒərs'isiz ðuŋ ʒ'udʒə
Una pinteta d'or sobre es coixí	unə piŋt'ətə ð'ɔ s'oβr əs kof'i
Sa seva germana no s'acaba de decidir	sə s'vevə zərm'anə n'o sək'abə ðə ðəsi'ði
En Joan diu que mos vàrem conèixer a ses festes des	ən ʒu'aŋ d'iw kə moz v'arəŋ kon'əʃ ə səs f'estəz ðəs



Provisió de dades

- Corpus textual
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - Campanya recollida de veu CM
 - Corpus de veu amb transcripció
 - Corpus de veu no alineat
 - Corpus de traducció parla a parla
- **Corpus paral·lels per TA**
 - **Corpus entrenament**
 - Corpus avaluació
- Estratègia de provisió de dades

Corpus paral·lels per entrenament de models de TA

Llengües	Num. de Frases	Origen de les dades	Disponible a
Català-Espanyol	85.953.513	Diverses fonts	-
Català-Anglès	22.592.622	Diverses fonts + AINA	projecte-aina/CA-EN_Parallel_Corpus
Català-Francès	18.634.844	OPUS	projecte-aina/CA-FR_Parallel_Corpus
Català-Portuguès	9.892.953	OPUS+Softcatalà (augmentat)	projecte-aina/CA-PT_Parallel_Corpus
Català-Italià	9.482.927	OPUS	projecte-aina/CA-IT_Parallel_Corpus
Català-Alemanys	9.530.709	OPUS +Softcatalà (augmentat)	projecte-aina/CA-ZH_Parallel_Corpus
Català-Xinès	6.833.114	OPUS (augmentat)	projecte-aina/CA-DE_Parallel_Corpus



Llengües	Num. de Frases	Origen de les dades	Disponible a
Català-Gallec	123.890.107	NOS + AINA	projecte-aina/CA-GL_Parallel_Corpus
Català-Euskera	96.692.186	GAITU + AINA	projecte-aina/CA-EU_Parallel_Corpus

En col·laboració amb
ILENIA
IMPULSO DE LAS LENGUAS EN LA INTELIGENCIA ARTIFICIAL



Provisió de dades

Corpus paral·lels per adaptació i avaluació de models de TA

- Corpus textual
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - Campanya recollida de veu CM
 - Corpus de veu amb transcripció
 - Corpus de veu no alineat
 - Corpus de traducció parla a parla
- **Corpus paral·lels per TA**
 - Corpus entrenament
 - **Corpus avaluació**
- Estratègia de provisió de dades

Corpus	Font	Llengües	Frases	Domini	Disponible a	Llicència
GEnCaTa	crawling .gencat	ca, en	38.595	Administratiu	ELRC-Share	CC0 4.0
Corpus bilingüe CA-EN de la CE	documents bilingües	ca, en	46.048	Administratiu	ELRC-Share	CC-BY 4.0
Col·lecció de corpus CA-EN de l'AP	documents bilingües	ca, en	36.116	Diversos	ELRC-Share	CC-BY 4.0
Col·lecció de corpus CA-ES de l'AP	documents bilingües	ca, es	63.773	Diversos	ELRC-Share	CC-BY 4.0
TaCon	document bilingüe	ca, es, gl, eu	1.314	Legal	ELRC-Share	CC-BY 4.0
Cyber MT test set	traducció	ca, en, es	1.715	Ciberseguretat	ELRC-Share	CC-BY-NC-SA 4.0
Catalan WMT2013	traducció	several	3.000	Notícies	ELRC-Share	CC-BY 4.0
Must-SHE Cat*	traducció	en, es, ca	1.046	Biaix de gènere		CC BY NC ND 4.0



*Permissos per publicar encara en tràmit



Provisió de dades

- Corpus textual
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - Campanya recollida de veu CM
 - Corpus de veu amb transcripció
 - Corpus de veu no alineat
- Corpus paral·lels per TA
- **Estratègia de provisió de dades**

Eines per al subministrament continu de dades de veu

- **Datapipe** una eina per facilitar l'adquisició dels continguts audiovisuals amb llicències obertes a la web. Projecte iniciat per la comunitat de programari lliure, que hem adaptat i millorat.

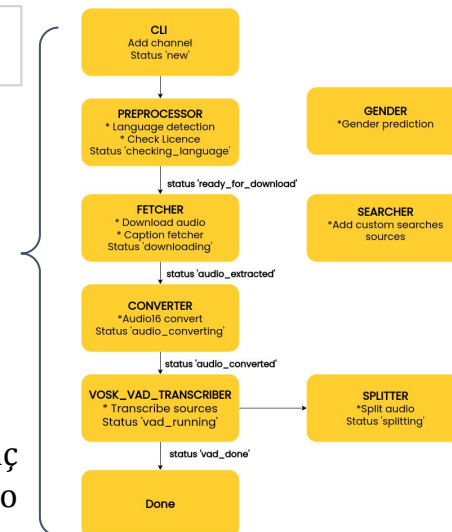
Accés: <https://github.com/projecte-aina/datapipe>

- **Found speech pipeline:** una eina per processar els continguts adquirits i generar un corpus de la parla alineat. El pipeline és capaç de generar dades per tasques de parla, a partir de transcripcions i/o subtítols.

Accés: TBA

L'objectiu és facilitar la **generació de datasets per a ASR, procesant continguts de forma automàtica**. Actualment, les estem usem en cinc contextos diferents:

- **Canals de youtube:** S'han identificat i descarregat continguts amb subtítols (~ 300 hores)
- **CCMA** (~ 4000 hores)
- **Parlament de Catalunya** (~ 4000 hores)
- **Corts Valencianes** (1200 hores)
- **IB3** (60 hores)



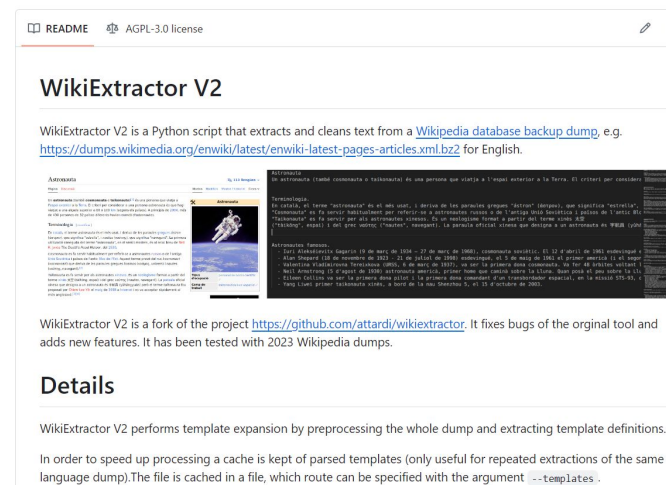


Provisió de dades

- Corpus textual
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - Campanya recollida de veu CM
 - Corpus de veu amb transcripció
 - Corpus de veu no alineat
- Corpus paral·lels per TA
- **Estratègia de provisió de dades**

Eines per al subministrament continu de dades de text

- Extractor de text de la **Viquipèdia** que interpreta correctament textos en una gran varietat de llengües, entre elles el català.
 - Stats: extracció de 689,141 documents amb més de 266M de paraules en català.
 - Accés: <https://github.com/langtech-bsc/Wikiextractor-V2/>



- Operacionalització de les dades obertes del **DOGC** mitjançant una pipeline automatitzada que fa ús de l'API de transparència de Catalunya.
 - Stats: fins a principis d'octubre, hem extret 30,369 publicacions en català amb 70M de paraules.
 - Accés: <https://github.com/projecte-aina/docg-pipeline>



Provisió de dades

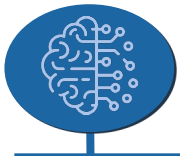
- Corpus textual
- Generació de dades anotades
 - Generació d'instruccions
- Corpus de veu
 - Campanya recollida de veu CM
 - Corpus de veu amb transcripció
 - Corpus de veu no alineat
- Corpus paral·lels per TA
- **Estratègia de provisió de dades**

Accions per a l'adhesió de noves dades

- Col·laboració amb el projecte [Parla'm](#) promogut per **Òmnium Cultural** per a la recollida de dades textuals i de veu.

Fins a la data les següents entitats han signat el *Compromís de cessió de continguts*

- Llenguaferrits (text i àudio)
- BIFIDUS PRODUCCIONS SL (text i àudio)
- Unió de Pagesos (text i àudio)
- Grup Món (text)
- Núvol (text)
- Betevé (àudio)



Models de la llengua



HUGGING FACE

<https://huggingface.co/projecte-aina>

Models Transformers

Generació de models de la llengua

- **Models Transformers de la llengua**
- Adaptació de models a tasques específiques (*fine tuning*)
- Participació en models multilingües i massius en col·laboració amb altres iniciatives



FLOR (Bloom) 760M

- Model trilingüe (ca,es,en) generatiu basat en Bloom entrenat amb **26B de tokens**.
- <https://huggingface.co/projecte-aina/FLOR-760M>



FLOR (Bloom) 1.3B

- Model trilingüe (ca,es,en) generatiu basat en Bloom entrenat amb **26B de tokens**.
- <https://huggingface.co/projecte-aina/FLOR-1.3B>



FLOR (Bloom) 6.3B

- Model trilingüe (ca,es,en) generatiu basat en Bloom entrenat amb **140B de tokens!!!**.
- <https://huggingface.co/projecte-aina/FLOR-6.3B>
- Sandbox: <https://huggingface.co/spaces/projecte-aina/flor-6.3b-inference>



Flor 6.3B Instruït

- Model Flor (Bloom) 6.3B pel català, anglès i castellà a partir de l'original
- Accés: <https://huggingface.co/projecte-aina/flor-6.3B-instruct>
- Sandbox: [HF Spaces](https://huggingface.co/spaces/projecte-aina/flor-6.3b-inference)



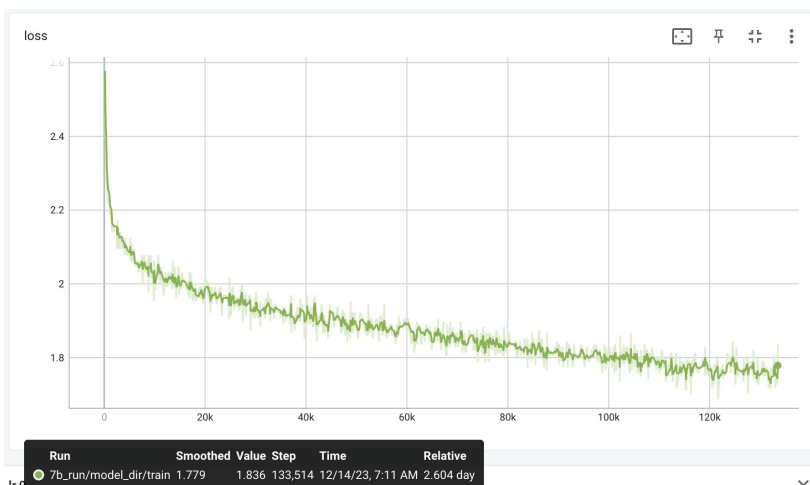
DeBerta_ca

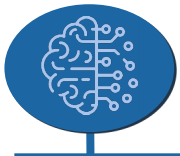
- Entrenat amb la segona versió del corpus textual català
- Accés: (en proves)



Sentence Encoders

- Entrenaments amb més dades del català dels embeddings
- https://huggingface.co/projecte-aina/ST-NLI-ca_paraphrase-multilingual-mpnet-base





Models de la llengua



HUGGING FACE

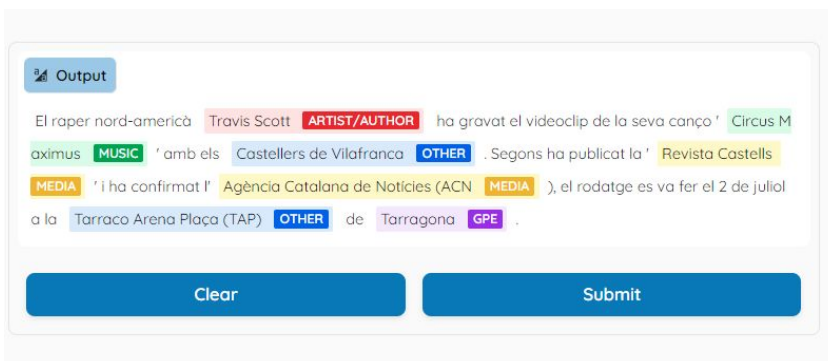
<https://huggingface.co/projecte-aina>

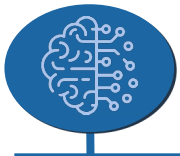
Models adaptats a tasques específiques

Generació de models de la llengua

- Models Transformers de la llengua
- **Adaptació de models a tasques específiques (*fine tuning*)**
- Participació en models multilingües i massius en col·laboració amb altres iniciatives

- **Catalan BERTa (RoBERTa-large) finetuned for Named Entity Recognition**
 - Tasca: Named Entity Recognition
 - Accés: https://huggingface.co/projecte-aina/multiner_ceil
- **RoBERTa-base v2 fine-tuned for QA**
 - Tasca: Pregunta/resposta (QA)
 - Accés: <https://huggingface.co/projecte-aina/roberta-base-ca-v2-cased-qa>
- **RoBERTa-base v2 fine-tuned for TE**
 - Tasca: Implicació textual
 - Accés: <https://huggingface.co/projecte-aina/roberta-base-ca-v2-cased-te>
- **RoBERTa-base v2 fine-tuned for STS**
 - Tasca: Similitud textual semàntica
 - Accés: <https://huggingface.co/projecte-aina/roberta-base-ca-v2-cased-sts>
- **RoBERTa-base v2 fine-tuned for Paraphrase Detection**
 - Tasca: Paràfrasi
 - accés: <https://huggingface.co/projecte-aina/roberta-large-ca-paraphrase>
- **Word embeddings Floret per al català, v1.0**
 - Accés: <https://zenodo.org/record/733033>



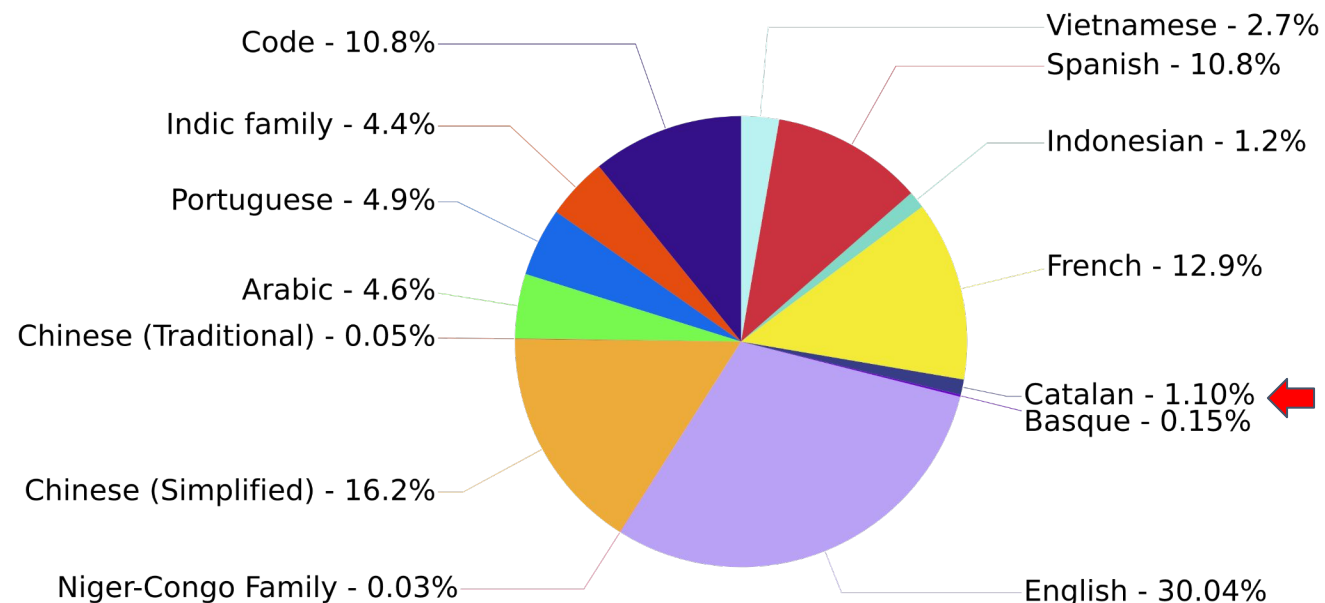


Models de la llengua

Generació de models de la llengua

- Models Transformers de la llengua
- Adaptació de models a tasques específiques (*fine tuning*)
- **Participació en models multilingües i massius en col·laboració amb altres iniciatives**

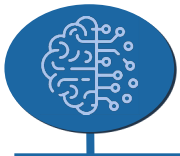
Liderant GPT3 massiu multilingüe



a BigScience initiative
BLOOM
176B params · 59 languages · Open-access

<https://huggingface.co/bigscience/bloom>

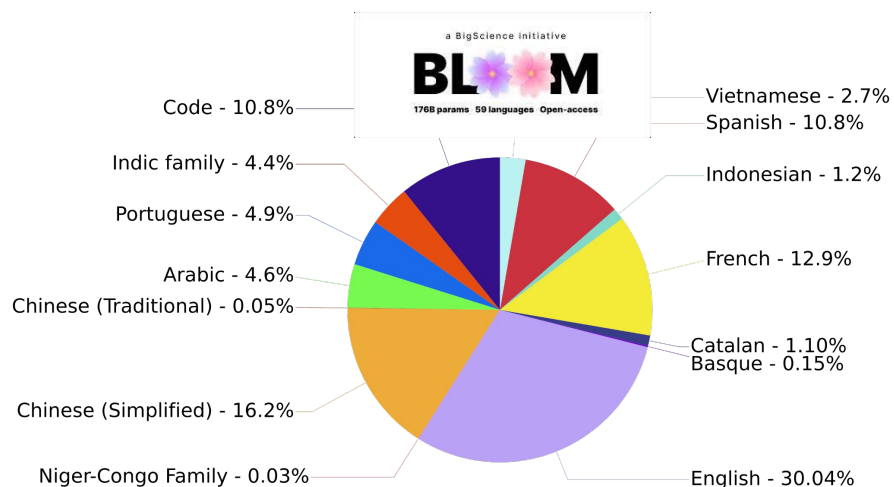
350B tokens (1,1% ca)



Models de la llengua

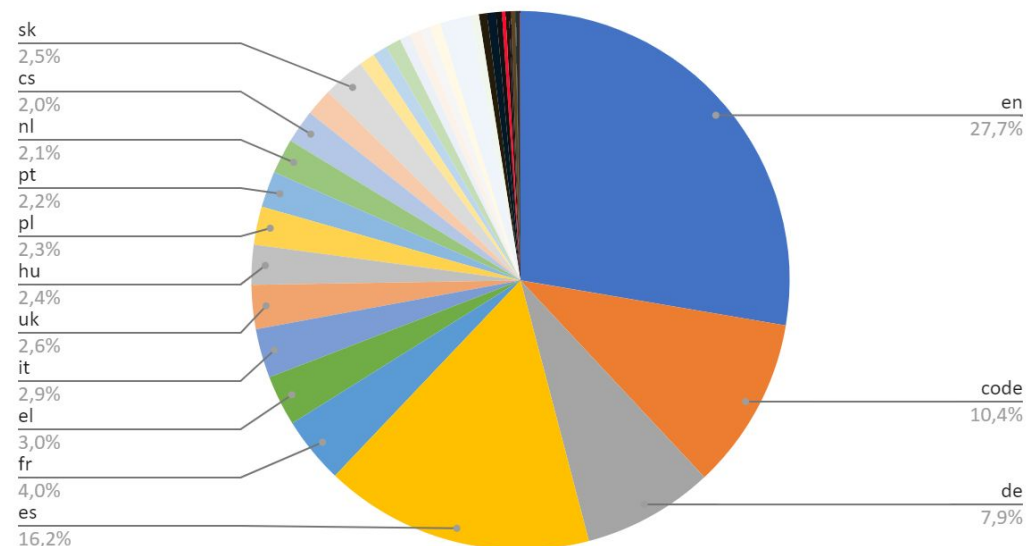
Generació de models de la llengua

- Models Transformers de la llengua
- Adaptació de models a tasques específiques (*fine tuning*)
- Participació en models multilingües i massius en col·laboració amb altres iniciatives



Liderant GPT3 massiu multilingüe

ca 1.86%



32 natural languages

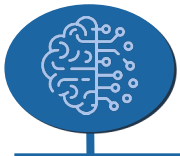
2,2T tokens (=2.200B) (1,86% ca)

<https://huggingface.co/bigscience/bloom>

46 natural languages

350B tokens (1,1% ca)





Models de la parla

Generació de models de la parla

- **Models de tecnologies de la parla**
- La presència dels models als ecosistemes i plataformes d'impacte

Motors de transcripció fonètica: Hem desenvolupat dos motors de transcripció fonètica, per convertir lletres en fonemes per 4 dialectes del català: central, nord-occidental, valencià i balear. Estan integrats als motors de TTS més utilitzats.

- <https://github.com/projecte-aina/espeak-ng/>
- <https://github.com/fedecosta/gruut>
en col·laboració amb la UPC

S'han utilitzat els enregistraments fets al **Common Voice** per crear dos models de la parla

Reconeixement de Parla - Nvidia Nemo: Alta precisió, mida mitjana, fàcilment desplegable amb el framework de **Nemo**.

- 36.5M paràmetres, WER of 6.684.
- <https://huggingface.co/projecte-aina/stt-ca-citrinet-512>

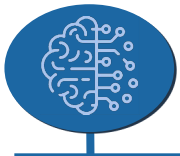


Síntesi de la Parla - Multiparlant: Model de síntesi de la parla d'alta qualitat i d'alt rendiment.

S'ha entrenat amb 255 veus i 487 hores d'enregistraments amb diverses variants dialectals.

A més de CV, aprofita els corpus de Festcat i de OpenSLR69. Fàcilment desplegable amb el framework de **Coqui**.

- <https://huggingface.co/projecte-aina/tts-ca-coqui-vits-multispeaker>
- integrat a la demo del bot d'AINA



Models de la parla

Generació de models de la parla

- Models de tecnologies de la parla
- **Presència dels models als ecosistemes i plataformes d'impacte**



Rasa is the leading platform for transforming how people interact with organizations through extensible conversational AI

<https://rasa.com/>



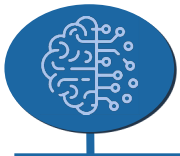
COQUI is a library for advanced Text-to-Speech generation

<https://github.com/coqui-ai/TTS>



NVIDIA NeMo, part of the NVIDIA AI platform, is a framework for building, training, and fine-tuning GPU-accelerated speech and natural language understanding (NLU)

<https://developer.nvidia.com/nvidia-nemo>



Models de Traducció

Generació de models de traducció

- Models de traducció
- Avaluació

En col·laboració



Traductor gallec → català

Accés: <https://huggingface.co/projecte-aina/mt-aina-gl-ca>

Traductor basc → català

Accés: <https://huggingface.co/projecte-aina/mt-aina-eu-ca>

- **Traductors català ↔ espanyol**
 - Accés (ca → es) : <https://huggingface.co/projecte-aina/mt-aina-ca-es>
 - Accés (es → ca) : <https://huggingface.co/projecte-aina/mt-aina-es-ca>
 - Accés (ca → es àmbit admin-legal) : <https://huggingface.co/projecte-aina/mt-aina-ca-es-adm>

- **Traductors català ↔ anglès**
 - Accés (ca → en) : <https://huggingface.co/projecte-aina/mt-aina-ca-en>
 - Accés (en → ca) : <https://huggingface.co/projecte-aina/mt-aina-en-ca>



Traductors català → portuguès

- Accés (ca → pt) : <https://huggingface.co/projecte-aina/mt-aina-ca-pt>
- Accés (pt → ca) : <https://huggingface.co/projecte-aina/mt-aina-pt-ca>



Traductors català ↔ italià

- Accés (ca → it) : <https://huggingface.co/projecte-aina/mt-aina-ca-it>
- Accés (it → ca) : <https://huggingface.co/projecte-aina/mt-aina-it-ca>



Traductors català ↔ francès

- Accés (ca → fr) : <https://huggingface.co/projecte-aina/mt-aina-ca-fr>
- Accés (fr → ca) : <https://huggingface.co/projecte-aina/mt-aina-fr-ca>



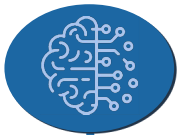
Traductors català ↔ alemany

- Accés (ca → de) : <https://huggingface.co/projecte-aina/mt-aina-ca-de>
- Accés (de → ca) : <https://huggingface.co/projecte-aina/mt-aina-de-ca>



Traductors català ↔ xinès

- Accés (ca → zh) : <https://huggingface.co/projecte-aina/mt-aina-ca-zh>
- Accés (zh → ca) : <https://huggingface.co/projecte-aina/mt-aina-zh-ca>

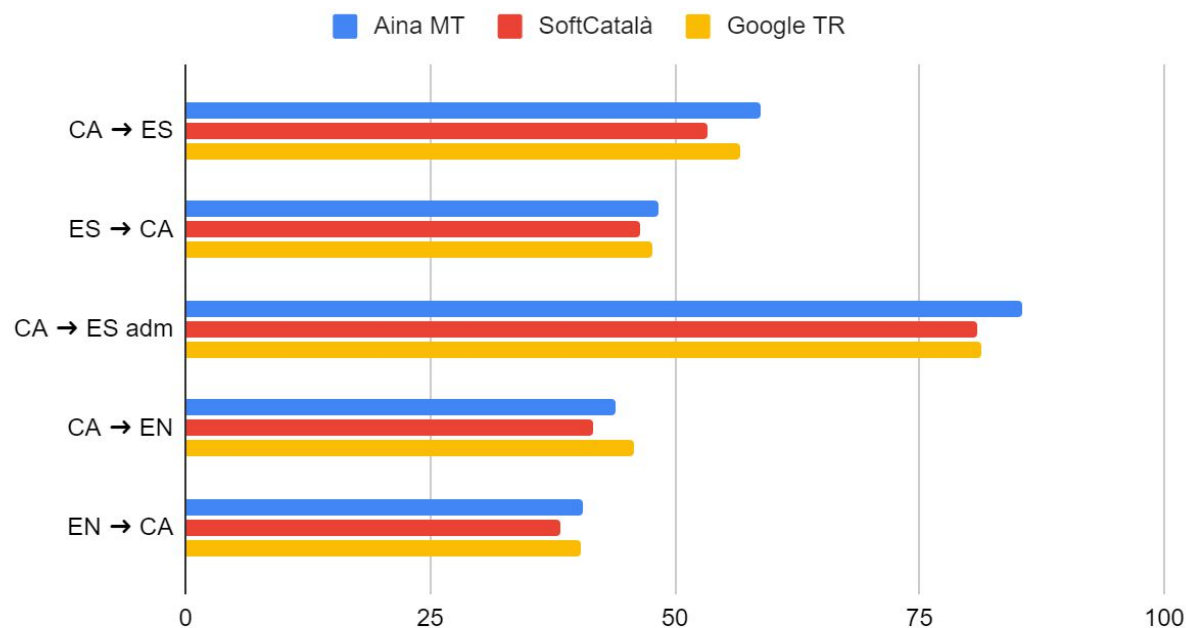


Models de Traducció

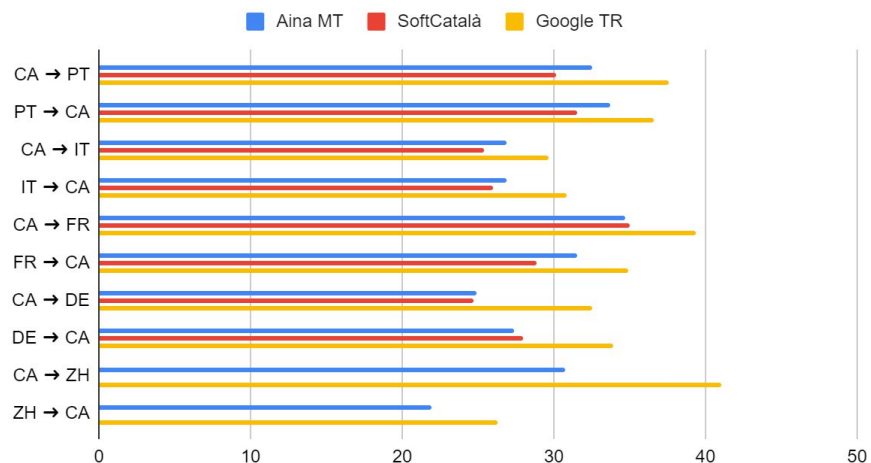
Generació de models de traducció

- Models de traducció
- **Avaluació**

Aina MT, SoftCatalà i Google TR



Aina MT, SoftCatalà i Google TR



Resultats de l'avaluació dels models català -anglès i català - castellà (mètrica BLEU, fent mitjana sobre diferents dominis textuais).

Resultats de l'avaluació de les versions inicials dels nous models (mètrica BLEU, fent mitjana sobre diferents dominis textuais).



Models de Traducció



The screenshot shows a web browser window with the URL 'e-aina/translator'. The page has a header with 'private', 'Running', and 'Logs' indicators, and navigation links for 'App', 'Files', and 'Community'. The main interface is divided into two columns. The left column is titled 'Source Language' and has a dropdown menu currently set to 'Catalan'. Below this is a large text input area with the placeholder text 'Enter a text here to translate.' and a character count '0 / 1000'. The right column is titled 'Target Language' and features radio buttons for 'Spanish', 'English', 'French', 'German', 'Italian', and 'Portuguese', with 'Spanish' selected. Below the target language options is another large text area. At the bottom of the interface are two blue buttons: 'Clear' and 'Submit'.



10 motors de traducció entre el català i diverses llengües europees, accessibles des del espai de Traducció del Projecte AINA a Hugging Face



Avaluació i benchmarking

- Plataforma de benchmarking
- Avaluació models generatius

Plataforma d'avaluació contínua de models amb avaluació extrínseca sobre 7 tasques diferents

<https://club.aina.bsc.es/>

Leaderboard

CLUB tests the ability of a system in the Catalan language. Below are the results of the different models.

Rank	Model	Submitted By	URL	Score	NER (F1)	POS (F1)	STS-ca (Comb.)	TeCla (Acc.)	TE-ca (Acc.)	CatalanQA (F1/EM)	XQuAD-ca (F1/EM)
1	RoBERTa-large-ca-v2	Projecte AINA		80.41	89.76	99.02	83.41	75.46	83.61	90.48/77.94	72.77/51.2
2	RoBERTa-base-ca-v2	Projecte AINA		79.29	89.27	98.95	79.07	74.26	83.14	89.37/75.64	72.79/51.1
3	mBERT	Projecte AINA		76.13	86.87	98.83	74.26	69.90	74.63	86.90/74.19	68.79/50.8
4	XLM-RoBERTa	Projecte AINA		71.23	86.31	98.89	61.61	70.14	33.30	88.17/75.93	72.55/54.1

Showing 1 to 4 of 4 entries

SEND YOUR RESULTS



Avaluació i benchmarking

- Plataforma de benchmarking
- Avaluació models generatius

```
python main.py \  
  --model hf-causal \  
  --model_args projecte-aina/Aguila-7b \  
  --tasks parafraseja, .....\  
  --device cuda:0
```

 [EleutherAI / lm-evaluation-harness](#)

Incorporació del català al nou LM-evaluation-Harness (EleutherAI)

- Incorporació de 15 datasets d'avaluació al LM-Harness
 - Bebebe (reading Comprehension),
 - Flores (TA)
 - CaBreu (resum)
 - CatalanQA, XQuad, CoQCat (QA)
 - PawsX, Parafraseja (Parafrasis)
 - XNLI, TeCa (NLI)
 - ARCChallenge i ARCEasy (raonament i coneixement sobre ciències).
 - PIQA, (raonament lògic sobre el món físic).
 - OpenBookQA (coneixements del món i certa capacitat de raonament).
 - MGSM (raonament matemàtic).
 - XStory Cloze (comprensió del llenguatge i raonament espacial i temporal propi de les narracions).
 - LMentry (avaluació capacitats lingüístiques)
- treballant en avaluació humana dels models generatius.

traduït



IA ètica

- Anàlisi de biaixos
- Lluita contra la desinformació

- **Identificació i filtratge de dades no desitjables** a les dades d'entrenament: contingut pornogràfic
 - integració a la Pipeline de preprocés
- **Inclusió string GUID** identificador únic per a poder excloure de futurs *crawlers* els datasets d'avaluació, garantint l'avaluació justa (Fair evaluation).
- Generació **corpus de factualitat** equivalent al TruthfulQA '**VeritasQA**' (360 preguntes/respostes en ca/en/es/gl/eu/. no conté coneixement 'localitzat', sense connexió geogràfica, 'traduïble', comprovant factualitat i amb un source de confiança)
- Avaluacions de biaix de gènere amb el nou corpus Must-SHE de TA



Serveis lingüístics

- **Accés i desplegaments**
- Integració de models en frameworks i plataformes de referència.

- El serveis dins **Spaces** en **Huggingface** tenen oberta **API** gratuïta, des de python, javascript o curl.
- Es poden fer servir per fer proves no-intensives des de codi.

Serveis lingüístics accessibles des de



Spaces 'Demos' & APIs

NOU **Model Flor 6.3 Instruct:** <https://huggingface.co/spaces/projecte-aina/flor-6.3b-instruct>

NOU **Model Flor 6.3:** <https://huggingface.co/projecte-aina/FLOR-6.3B>

NOU **Transcripció fonètica** <https://huggingface.co/spaces/projecte-aina/transcripcio-fonetica-catala>

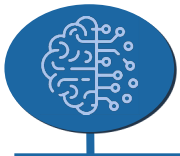
api_name: `/submit_input`

```
from gradio_client import Client

client = Client("https://projecte-aina-flor-6-3b-inference.hf.space/--replicas/akqui/")
result = client.predict(
    "Hello!!", # str in 'Input' Textbox component
    1, # float (numeric value between 1 and 200) in 'Max tokens' Slider component
    0.1, # float (numeric value between 0.1 and 10) in 'Repetition penalty' Slider component
    1, # float (numeric value between 1 and 100) in 'Top k' Slider component
    0.01, # float (numeric value between 0.01 and 0.99) in 'Top p' Slider component
    True, # bool in 'Do sample' Checkbox component
    1, # float (numeric value between 1 and 8) in 'Beams' Slider component
    0, # float (numeric value between 0 and 1) in 'Temperature' Slider component
    api_name="/submit_input"
)
print(result)
```

• Return Type(s)

```
# str representing output in 'Output' Textbox component
```

Serveis lingüístics

- Accés i desplegaments
- Integració de models en frameworks i plataformes de referència.

spaCy

[Spacy 3.7](#) permet ara la integració directa dels models Transformers d'AINA, inclosos LLMs: [spacy-huggingface-pipelines](#) | [spacy-transformers](#) | [spacy-llm](#)

```
import spacy
nlp = spacy.load("ca_core_news_trf")

nlp.add_pipe("hf_token_pipe", config={"model": "projecte-aina/multiner_ceil"})
doc = nlp("El Barça ha fitxat a Messi pel 2024.")
for e in doc.ents: print(e,e.label_)
#Barça organization-sportsteam
#Messi person-athlete
```

aws

S'han preparat **Notebooks** per fer servir al còmput al núvol de AWS els recursos d'AINA per fer **inferència** i **fine-tuning**. Accés [sagemaker-examples](#)

- Inference: [Aguila 7b Hugging Face Large Model Inference - TGI](#) shows how to deploy common large language models such as projecte-aina/aguila-7b, using Hugging Face Text Generation Inference (TGI) Deep Learning Container with SageMaker
- Fine-tuning: [Aguila 7b fine-tuning with instruction dataset](#) shows how to fine-tune the falcon 7B aguila model projecte-aina/aguila-7b, using an instructional dataset (in this case an example from the InstructCat collection) with a EC instance with Sagemaker.

- **Anonimitzador** de continguts generats per usuaris * <https://github.com/TeMU-BSC/AnonymizationPipeline>
- **TTS**: RestFUL API and web interface to serve coqui TTS models <https://hub.docker.com/r/projecteaina/tts-api> <https://github.com/projecte-aina/tts-api>
- **TA**: RestFUL API for serving machine translation models <https://github.com/projecte-aina/mt-api>
- **RASA**: Exemple d'assistent conversacional <https://github.com/projecte-aina/minibot>



* Avaluació en col·laboració amb **1million bot**



On trobar-ho?

<https://bit.ly/AINAKit>

The screenshot shows the Aina Kit website home page. At the top left, it says "Aina Kit - Home". On the right, there is a search bar and a "Try Notion" button. The main header features the "Aina" logo in red and white, with a house icon below it. To the right of the logo are the logos for "Barcelona Supercomputing Center" and "Generalitat de Catalunya", with the text "Finançat per" above the latter. Below the header is a navigation menu with icons and labels for HOME, MODELS, DATASETS, TEST, DEMOS, and PARTNERS. The main content area starts with the heading "Aina Kit - Home" and a paragraph: "Aina Kit és el recull de documentació i recursos que el projecte Aina posa a disposició de les persones desenvolupadores de productes i serveis de d'intel·ligència artificial (IA) i tecnologies del llenguatge (TL) en català".



On trobar-ho?

<https://bit.ly/AINAKit>



GitHub

<https://github.com/projecte-aina>



<https://huggingface.co/projecte-aina>

zenodo

<https://zenodo.org/communities/catalan-ai>



<https://elrc-share.eu/repository/search/>



docker

<https://hub.docker.com/u/projecteaina>

Aina Kit - Home

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

Finançat per
Generalitat de Catalunya

HOME MODELS DATASETS TEST DEMOS PARTNERS

L'Aina Kit

Aina Kit és el recull de documentació i recursos que el projecte Aina posa a disposició de les persones desenvolupadores de productes i serveis de d'intel·ligència artificial (IA) i tecnologies del llenguatge (TL) en català.



On trobar-ho?



<https://github.com/projecte-aina>



HUGGING FACE

<https://huggingface.co/projecte-aina>



<https://zenodo.org/communities/catalan-ai>



<https://elrc-share.eu/repository/search/>



<https://hub.docker.com/u/projecteaina>

Collections 6

TEXT_Models

Encoders / Decoders models, foundational, pretrained or fine-t...

projecte-aina/aguila-7b
Text Generation • Updated Oct 31 • 1.1k • 31

projecte-aina/roberta-base-ca-v2
Fill-Mask • Updated Dec 23, 2022 • 10 • 2

projecte-aina/longformer-base-4096-ca-v2
Fill-Mask • Updated Dec 15, 2022 • 4

projecte-aina/roberta-large-ca-v2

MT_Models

Machine Translation models

projecte-aina/mt-aina-ca-es-adm
Updated 8 days ago

MT_Datasets

Machine Translation datasets

projecte-aina/CA-PT_Parallel_Corpus
Updated 9 days ago • 2

projecte-aina/ca_zh_wikipedia
Viewer • Updated Jan 9 • 34 • 3



HUGGING FACE

^ Collapse

TEXT_Datasets

Datasets for fine-tuning, instruction and evaluation of text mo...

projecte-aina/catalanqa
Updated 13 days ago • 174 • 1

projecte-aina/viquiquad
Viewer • Updated Sep 13 • 88 • 1

projecte-aina/ancora-ca-ner
Viewer • Updated Sep 13 • 69

projecte-aina/CaWikiITC

SPEECH_Datasets

Speech models and datasets

projecte-aina/parlament_parla
Preview • Updated Sep 13 • 44 • 1

SPEECH_Models

ASR, TTS and other speech/audio related tasks

projecte-aina/tts-ca-coqui-vits-multispea...
Updated Dec 19, 2022 • 3

projecte-aina/stt-ca-citriNET-512
Automatic Speech Recognition • Updated D... • 16 • 1



On trobar-ho?



<https://github.com/projecte-aina>



HUGGING FACE

<https://huggingface.co/projecte-aina>



<https://zenodo.org/communities/catalan-ai>



<https://elrc-share.eu/repository/search/>



<https://hub.docker.com/u/projecteaina>



HUGGING FACE

- Train
- Deploy
- Use in Transformers

Downloads last month
1,100

- AutoTrain
Fine-tune this model without code
- Amazon SageMaker
Use SageMaker for optimized training



Safetensors Model size 6.85B params Tensor type F32

- Train
- Deploy
- Use in Transformers

Downloads last month
122

Inference API

Token Classification

Your sentence here...

Compute

- Inference API
Free API for fast prototyping
- Inference Endpoints
Production-ready API deployments
- Amazon SageMaker
Optimized deployments with SageMaker
- Spaces
Deploy as a Gradio app in one click
- Azure ML
Optimized deployments with AzureML



Comunicació



Ja som més de 3.100!

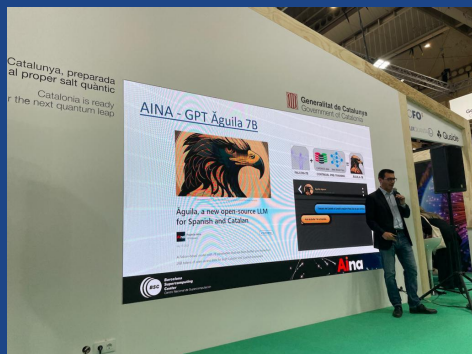
@projecte_aina



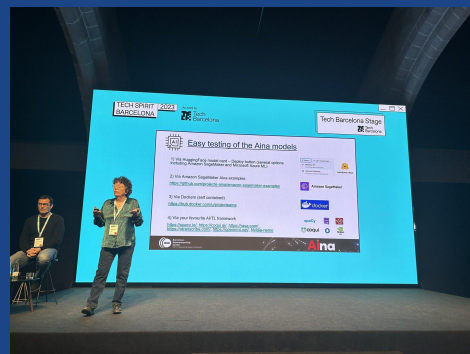
+ 240 seguidors

Projecte Aina

ESDEVENIMENTS



Presentació al Smart City Expo World Congress 2023



Presentació al Tech Spirit Barcelona 2023



WEB ACTUALITZADA, JA DISPONIBLE!

<https://projecteaina.cat/tech/>



Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

Aina



Prototips i demostradors

<https://aina.bsc.es/>

Aina

Bot

chatbot speech voice

Demostració d'incorporació de funcionalitats de veu a un xatbot.

Spacy

text classification similarity tokenization

Demostrador de les capacitats de les cadenes de processament del llenguatge natural i models Spacy implementats dins del Projecte AINA.

ViquiQA

question answering wikipedia catalan

Demostrador del model de Pregunta i Resposta entrenat amb el dataset CatalanQA, fent servir la Viquipèdia en català.

Traductor

machine translation catalan english

Traductors automàtics entre català i castellà (text general i d'especialitat administratiu-legal) i entre català i anglès (text general).

oTranscribe+

speech recognition transcription catalan

Aplicació web amb reconeixement de la parla gratuïta i privada per a transcriure entrevistes enregistrades.

CLUB

model benchmark catalan

Plataforma d'avaluació comparativa de models de llengua per al català.

TTS

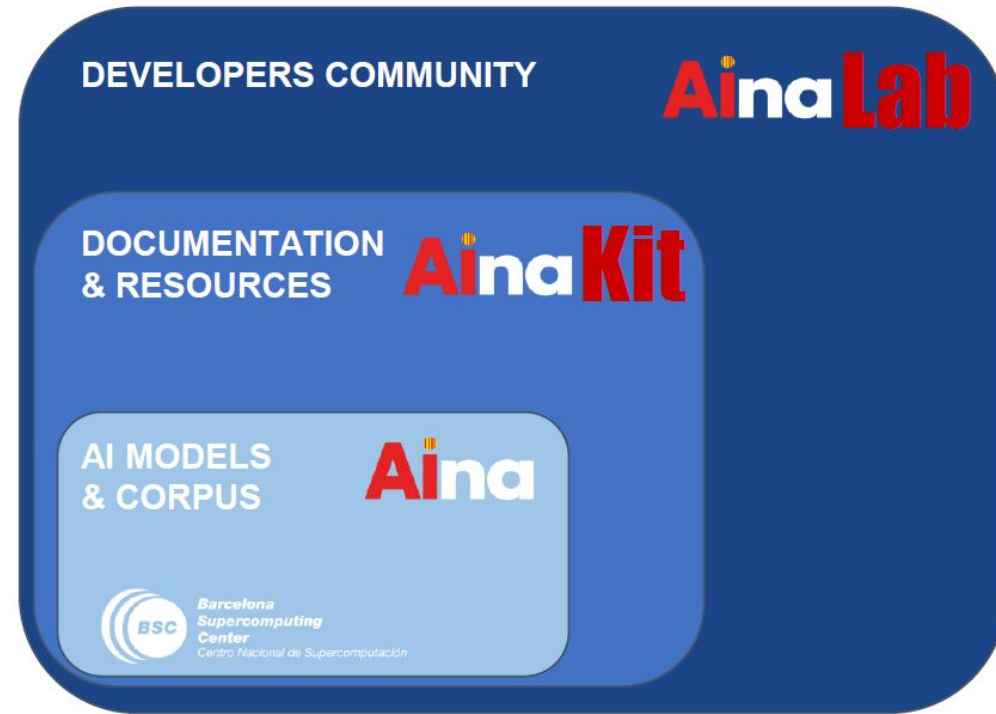
text-to-speech catalan

Demostrador del motor de síntesi de la parla multi parlant.

Aina Open Innovation Contest 2024

Accelerating the ecosystem on 2024

- Aina “Open Innovation Contest” (Q1 2024) will set in motion the developers community around the available open resources
 - ✓ Technical validation
 - ✓ Social impact evidence
 - ✓ Ecosystem development



Aina Open Innovation Contest 2024

2) Dedicated site for “expressions of interest” for open innovation contest

- https://twitter.com/projecte_aina
- <https://www.linkedin.com/company/projecte-aina/>

STAY TUNED!



Prototips i demostradors

Models Fundacionals i Instruïts

Demostració dels models generatius FLOR basats en BLOOM, tant el model fundacional com el prototip instruït, amb la possibilitat de modificar els paràmetres de generació.

Model Fundacional:

<https://huggingface.co/spaces/projecte-aina/flor-6.3b-inference>

Model Instruït:

<https://huggingface.co/spaces/projecte-aina/flor6.3b-instruct>

Model Instruït Lite (1.3b, corre sense GPUs):

<https://huggingface.co/spaces/projecte-aina/flor1.3b-instruct>



Flor-6.3B Instruct (experimental)

✿ **Flor** is a 6.3B parameters multilingual LLM that has been trained on a massive mixture of Spanish, Catalan and English data. It is a new open-source Large Language Model (LLM), licensed for both research and commercial use. It uses the [Bloom-7b](#) model as a starting point, a state-of-the-art multilingual language model.

✔ **Intended use:** This is the fine-tuned version of the model, trained with [InstruCatPlus](#) a merge of instruction datasets in English, Spanish and Catalan: It follows instructions for question answering, creative writing, paraphrase creation, named-entity detection, summarization, etc., but has not been refined yet with a Reinforcement Learning from Human Feedback (RLHF) process to align it with human behaviors or preferences. It is NOT a conversational chatbot, and won't keep a running memory of the interactions.

You can **ask questions or requests using the prompt**, and provide context if needed (such as text you want to summarize).

✂ **Instructional data:** [InstruCat](#) (Catalan), 216k instruction for tasks such as summarization, entailment, phrase generation, NER, etc. [Mentor_ES](#) (Spanish), 10k+ instructions in Spanish following the Dolly template, commissioned by the Language Technologies Group at BSC. [dolly](#) (English), a 11k subset of the Dolly dataset.

⚠ **Limitations:** This version is for beta testing only. The content generated by these models is unsupervised and might be judged as inappropriate or offensive. Please bear this in mind when exploring this resource.

•• **Learn more about Flor:** [HF official model card](#) and the [Instructed version](#).

Prompt	Context	Response
Qui es Raimon?	This field is optional.	Un cantautor valencià.



Prototips i demostradors



PoC: Assistent d'Informació Meteorològica (amb la col·laboració de la CCMA)

Demostració de generació automàtica de prediccions meteorològiques locals amb un assistent interactiu que fa servir una versió lleugera del model Flor1.3, amb col·laboració amb els meteoròlegs de la CCMA.


<https://huggingface.co/spaces/projecte-aina/BlooMeteo>

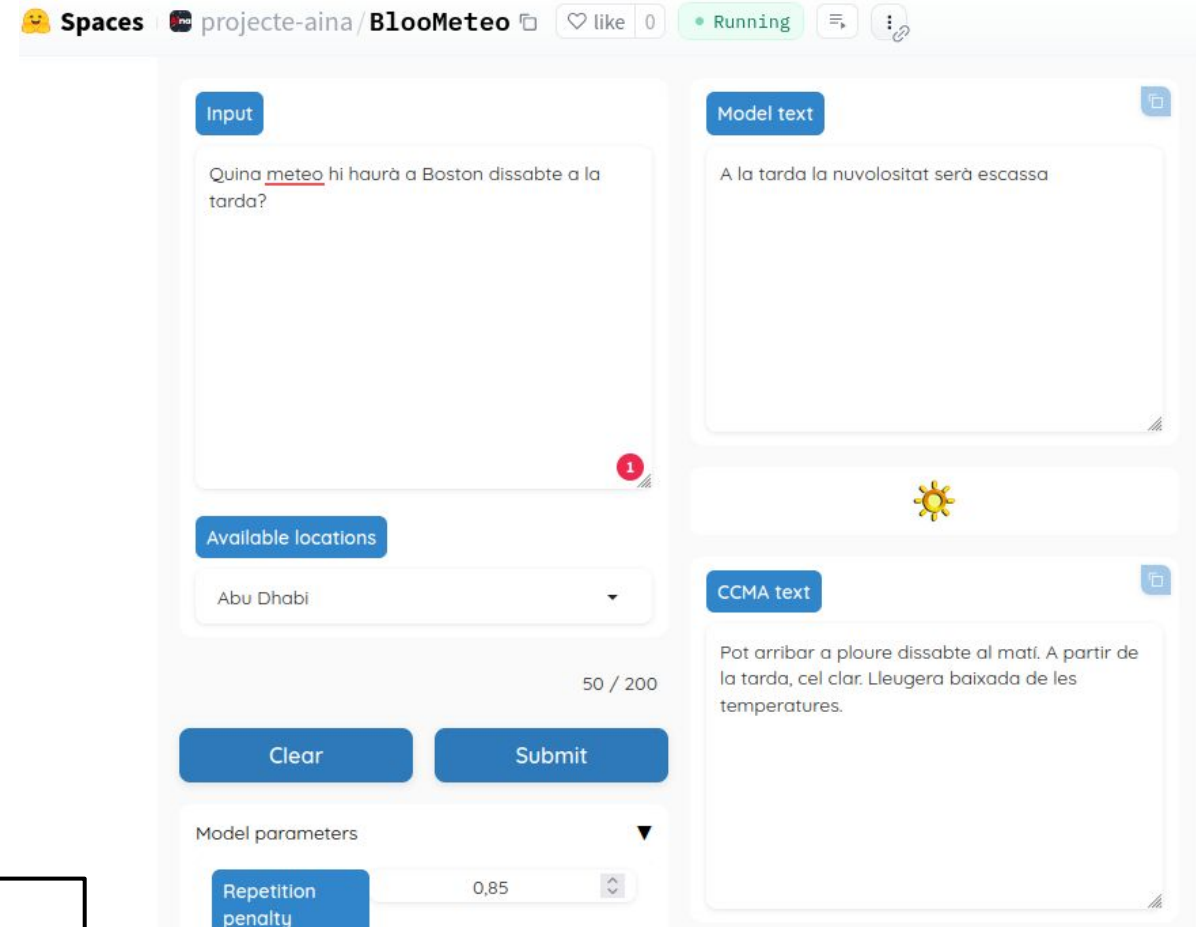
DeBERTA MultiNER (Multilingüe):

Demostració d'Anotació i Classificació d'Entitats en contextos multilingües, amb models massius DeBERTA, i el dataset CEIL.



https://huggingface.co/spaces/projecte-aina/multiner_demo

 Aviat, demo de **RAG (Retrieval-Augmented Generation)** per fer cerques d'informació en repositoris massius de documents fent preguntes en llenguatge natural



The screenshot shows the Hugging Face Spaces interface for the 'BlooMeteo' demo. At the top, it displays 'Spaces projecte-aina/BlooMeteo' with a 'like' count of 0 and a 'Running' status. The interface is divided into several sections:

- Input:** A text area containing the question 'Quina meteo hi haurà a Boston dissabte a la tarda?'. A red notification bubble with the number '1' is visible in the bottom right corner of this section.
- Available locations:** A dropdown menu currently showing 'Abu Dhabi'.
- Model text:** A text area displaying the model's output: 'A la tarda la nuvolositat serà escassa'.
- CCMA text:** A text area displaying additional information: 'Pot arribar a ploure dissabte al matí. A partir de la tarda, cel clar. Lleugera baixada de les temperatures.'

At the bottom, there are 'Clear' and 'Submit' buttons, and a 'Model parameters' section with a 'Repetition penalty' slider set to 0,85.



Prototips i demostradors



Síntesi de la parla en català (TTS)

El nostre nou model de TTS està considerablement millorat sobretot per la fiabilitat de pronunciació, i la qualitat d'àudio. Estem col·laborant amb Softcatalà i la CCMA per millorar. El nou model està integrat al nostre software tts-api, per facilitar les integracions a la producció. El software facilita l'ús dels models amb un millor rendiment.

- TTS-API: <https://hub.docker.com/r/projecteaina/tts-api>
- Demo TTS: <https://huggingface.co/spaces/projecte-aina/tts-ca-coqui-vits-multispeaker>
- Demo conversacional: <https://bot.aina.bsc.es/>



Identificació dels parlants

Vam desenvolupar un software que processa un àudio llarg per identificar qui està parlant quan. Utilitza enregistraments de referència per detectar la identitat de la persona, i relacionar-ho amb les pautes en que parla. El software està dissenyat per l'ús de l'administració pública. Aviat publicarem una versió genèrica del codi per l'ús de tothom.

POST /api/tts Tts

Text-to-Speech API endpoint.

This endpoint receives a TTSRequestModel object containing the voice and text to be synthesized. It performs the necessary processing to generate the corresponding speech audio and streams it back as a WAV audio file.

Parameters:

- request: TTSRequestModel - An object containing the voice and text data for synthesis.

Returns:

- StreamingResponse: A streaming response object that contains the synthesized speech audio as a WAV file.

Raises:

- SpeakerException: If the specified speaker ID is invalid.
- LanguageException: If the specified language is not supported.

Parameters Try it out

No parameters

Request body required application/json

Example Value | Schema

```
{
  "language": "ca-es",
  "voice": "string",
  "type": "string",
  "text": "string"
}
```



L'horitzó d'Aina: estatus del projecte

Actualització sobre el desenvolupament dels models lingüístics, recursos i objectius d'AINA

19 de desembre del 2023

