

# Norma, ús i interferència: biaixos lingüístics en els models de llenguatge en català

## Norm, use and interference: linguistic biases in Catalan language models

Mireia Almena Rodríguez ✉

Departament de Traducció i Ciències del Llenguatge, Universitat Pompeu Fabra

Thomas Brochhagen ✉ 

Departament de Traducció i Ciències del Llenguatge, Universitat Pompeu Fabra

### Resum

---

Els Models de Llenguatge Extensos influeixen cada cop més en la comunicació escrita, fet que planteja reptes per a llengües minoritàries com el català. Aquest estudi quantifica els biaixos lingüístics de sis models causals de llenguatge en català, analitzant-ne les preferències per construccions gramaticals normatives enfront de les no normatives, especialment sota la influència potencial de la interferència del castellà. A partir d'un corpus propi de parelles mínimes, hem avaluat models tant monolingües com multilingües comparant les seves preferències per a cada variant (no) normativa. Els resultats indiquen que no hi ha una diferència entre els models monolingües i multilingües en la seva preferència per les construccions normatives. Tanmateix, la possible interferència del castellà redueix de manera notable la preferència per les formes normatives en tots els models analitzats. Aquesta reducció és més pronunciada en models multilingües, indicant més influència del castellà. Aquests resultats suggereixen que els biaixos dels models reflecteixen la prevalença d'usos no normatius en les seves dades d'entrenament, a causa d'influència del castellà. Això subratlla la importància d'avaluar aquestes tecnologies per a informar la política lingüística i comprendre'n l'impacte en l'evolució de la llengua.

### Paraules clau

---

models de llenguatge; biaixos; norma; us; interferència; català

### Abstract

---

Large Language Models are increasingly influencing written communication. This poses challenges for minority languages such as Catalan. This study quantifies the linguistic biases of six Catalan language models, analyzing their preferences for normative versus non-normative grammatical constructions, especially for cases where there can be interference from Spanish. Using a corpus of minimal pairs, we evaluate both monolingual and multilingual models by comparing their preferences for each (non-)normative variant. The results indicate that there is no differ-

ence between monolingual and multilingual models in their preference for normative constructions. However, cases where there can be interference from Spanish markedly reduce the preference for normative forms across all analyzed models. This reduction is more pronounced in multilingual models, indicating stronger influence from Spanish. These findings suggest that the models' biases reflect the prevalence of non-normative usage in their training data, due to influence from Spanish. This underscores the importance of evaluating these technologies to inform language policy and understand their impact on language evolution.

### Keywords

---

language models; biases; norm; use; interference; catalan

## 1. Introducció

---

Les tecnologies del llenguatge tenen un paper cada vegada més important en la nostra vida quotidiana i en la manera com ens comuniquem. En particular, en pocs anys, s'ha estès l'ús de Models de Llenguatge Extensos generatius (LLMs, per les seves sigles en anglès). Ara per ara ja són eines que es fan servir per realitzar tasques rudimentàries, com redactar correus electrònics, o resumir i reestructurar textos. Podem suposar que, si l'adopció d'aquesta tecnologia segueix les tendències actuals, els LLMs formaran una part inextricable del pre- i post-processament del nostre ús de llenguatge escrit, com ho són ara els correctors automàtics d'ortografia i funcions d'emplenament automàtic.

Aquest desenvolupament té conseqüències per als llenguatges: si la proporció de text generada, o almenys modificada, per LLMs és cada vegada més gran (Liang et al., 2025), l'ús del llenguatge dels models —de biaixos a idiosincràsies estilístiques— no només retroalimentarà noves ge-

neracions de models, sinó que també pot afectar l'ús dels parlants. En altres paraules, els models no només reproduïen el llenguatge, sinó que també l'influencien. Això pot afectar més les llengües minoritàries, com el català, que són nodrides per un volum menor de contingut escrit —amb una esfera de parlants més petita— que altres llengües. En aquest estudi abordem aquesta temàtica, analitzant les preferències lingüístiques de models de llenguatge en català. D'una banda, volem quantificar els biaixos i les preferències per usos del llenguatge (no) normatius que tenen els models actuals. De l'altra, amb aquesta quantificació volem obrir un debat informat sobre el tipus de preferències lingüístiques que haurien de reproduir.

Un altre factor que pot afectar especialment les llengües minoritàries és l'ús generalitzat de models multilingües (Schut et al., 2025). Aquests models sovint s'han entrenat amb un volum més gran de dades —i sovint de més qualitat— en llengües com l'anglès, en comparació amb llengües com el català o el gallec. Aquest és el cas de tots els models multilingües catalans que analitzem aquí (vegeu §2). En conseqüència, les estructures gramaticals i el lèxic, entre altres, de les llengües amb més recursos poden afectar les llengües amb menys recursos. Diversos estudis ja han analitzat els biaixos lingüístics que presenten alguns models multilingües, especialment la seva tendència a afavorir estructures gramaticals i conceptuals angleses quan generen text en altres llengües (Papadimitriou et al., 2023; Wendler et al., 2024; Brinkmann et al., 2025; Schut et al., 2025). Però encara no sabem gaire sobre biaixos i interferències causades en aquests models entre altres llengües, com ho són el castellà i el català. En el cas del català, no queda clar fins a quin punt els models de llenguatge, tant monolingües com multilingües, fan servir estructures normatives, en comparació amb no-normatives. Amb aquest estudi, pretenem tancar aquesta bretxa.

Donat el fet que aquestes tecnologies acabaran influïnt en la nostra parla, la qual, al seu torn, alimentarà els nous models, considerem essencial rastrejar i avaluar les preferències lingüístiques dels models actuals; i analitzar fins a quin punt aquestes preferències estan influenciades per altres llengües. Això és el que ens proposem en aquest article: elaborem un corpus d'avaluació de parelles mínimes (no) normatives del català amb el qual avaluem sis models, tant mono- com multilingües.

## 2. Material i Mètodes

Tot el codi i el conjunt de dades que creem per avaluar models són disponibles a GitHub<sup>1</sup>. Comencem amb una descripció dels models triats i després descrivim les dades per avaluar-los.

**Models** Hem seleccionat un conjunt de sis models, dos de monolingües i quatre de multilingües, capaços de generar i predir contingut en català. Amb models *multilingües* catalans ens referim a models que indiquen oficialment haver estat entrenats, entre d'altres, amb dades en català. Els models multilingües que estudiem són: Salamandra-7B<sup>2</sup> (Gonzalez-Agirre et al., 2025), un model base entrenat en 35 llengües europees (39% anglès, 16% castellà, 2% català); BLOOM-7.1B<sup>3</sup> (Le Scao et al., 2023) un model base entrenat en 45 llengües (31% anglès, 11% castellà, 1% català); XGLM-7.5B<sup>4</sup> (Lin et al., 2022), un model base entrenat en 31 llengües (49% anglès, 5% castellà, 0.4% català); i ALIA-40B<sup>5</sup> (Gonzalez-Agirre et al., 2025), una variant més gran del model Salamandra, amb la mateixa cobertura lingüística.

Amb models *monolingües* catalans ens referim a models que han passat per un procés final d'aprenentatge que els ajusta específicament al català. Els models monolingües que estudiem són: CataLlama-8B v0.1<sup>6</sup>, basat en Llama-3 i entrenat durant un epoch amb 331 milions de tokens en català; i Aitana-6.3B<sup>7</sup>, un model basat en FLOR-6.3B i entrenat durant 2 epochs amb 1.3 milions de tokens en valencià. Excepte ALIA-40B i Salamandra-7B, tots aquests models han estat recentment avaluats pel que fa a les seves capacitats generals a IberoBench (Baucells et al., 2025). Però cap d'ells ha estat estudiat pel que fa a preferències de llenguatge (no) normatiu o possibles interferències.

La distinció que fem entre models monolingües i multilingües és habitual a la literatura. Cal esmentar, tanmateix, que la majoria dels models monolingües actuals han estat entrenats amb dades d'altres llengües, ja sigui en etapes inicials (abans del *fine tuning* per a la llengua objectiu), o bé de manera no intencionada, atès que les dades d'entrenament no s'han depurat fins al

<sup>1</sup><https://github.com/mireia-almena/NormaUsiInterferencia>

<sup>2</sup><https://huggingface.co/BSC-LT/salamandra-7b>

<sup>3</sup><https://huggingface.co/bigscience/bloom-7b1>

<sup>4</sup><https://huggingface.co/facebook/xglm-7.5B>

<sup>5</sup><https://huggingface.co/BSC-LT/ALIA-40b>

<sup>6</sup><https://huggingface.co/catallama/CataLlama-v0.1-Base>

<sup>7</sup><https://huggingface.co/gplsi/Aitana-6.3B>

punt de poder garantir que no s’hagin emprat dades d’una altra llengua en l’entrenament. Aquest és un fet a tenir en compte, particularment en avaluar models de llengües amb menys recursos. En conseqüència, i fet rellevant per al cas que ens ocupa, no és improbable que els models monolingües catalans hagin estat entrenats també, si més no parcialment, amb dades del castellà. Tornem a abordar aquesta qüestió en la discussió dels resultats.

Vam descartar dos altres models potencialment rellevants de l’avaluació: RoBERTa-ca i CataLlama-v0.2-Base. RoBERTa-ca<sup>8</sup> és un masked language model, amb un objectiu d’entrenament diferent del de la resta de models causals que avaluem. Atès que la nostra mètrica d’avaluació es basa en la negative log-likelihood (NLL, vegeu més avall), vam descartar models no causals perquè la manera de calcular la NLL és diferent, i aquesta diferència podria introduir soroll a l’anàlisi. CataLlama-v0.2<sup>9</sup>, una fusió entre CataLlama-v0.1 i Meta-Llama-3-8B-Instruct, va ser descartat perquè vam considerar que la fusió podria fer el model “més multilingüe” que la resta de membres de la categoria monolingüe, preferint avaluar la seva versió 0.1 sense fusió amb un model d’instrucció en anglès. Una anàlisi post hoc revela que aquestes dues decisions de descart no afecten qualitativament els nostres resultats: incloue RoBERTa-ca i CataLlama-v0.2 —analitzant 4 models multilingües i 4 de monolingües— suggereix les mateixes tendències que els resultats obtinguts amb l’exclusió d’aquests dos models (vegeu l’apèndix). Per consegüent, els resultats reportats a la secció §3 són robustos davant aquesta decisió.

**Corpus d’avaluació** Hem construït un corpus d’avaluació amb vint oracions per cada un de vuit tipus d’estructures gramaticals en català que poden donar lloc a la producció de respostes no normatives. Quatre d’aquestes estructures són susceptibles a generar ús fora de la norma per interferència amb el castellà. Les altres quatre, no relacionades amb la interferència, s’han inclòs per observar si alguns models tracten de manera diferent les construccions influïdes per una altra llengua en comparació amb aquells que són simplement desviacions de la norma. Il·lustrem cada tipus d’estructura amb una oració exemplar, indicant com a primera opció la normativa i com a segona la no normativa:

1. Em refereixo  $\{, a\}$  que va cantar.
2. He vist  $\{el, al\}$  cotxe nou.
3. Al restaurant fa molta olor  $\{de, a\}$  fregit.
4. Insisteixen  $\{a, en\}$  tornar cap a casa.
5. No sóc gens propens  $\{a, de\}$  enfadar-me per tonteries.
6. Van cancel·lar el vol  $\{amb, en\}$  motiu del mal temps.
7. Cal confiar  $\{en, amb\}$  el criteri dels experts.
8. Si no vas  $\{amb, en\}$  compte, et pots fer mal.

Aquest corpus permet avaluar la qualitat lingüística dels models a partir de frases en què coexisteix una forma normativa i una altra de no normativa, sovint present en l’ús dels parlants. Com il·lustrat dalt, totes les estructures incloses en els conjunts de dades es basen en construccions gramaticals relacionades amb preposicions, ja que es tracta d’una categoria tancada que ofereix menys variabilitat en les respostes dels models. Això facilita la comparació de les respostes com a normatives o no normatives.

Per a cadascuna de les 160 oracions, calculem la negative log-likelihood (NLL) assignada a la seva versió (a) normativa i (b) no normativa per cada model. En altres paraules, per comparar dos oracions, calculem el logaritme negatiu de les seves probabilitats:

$$\begin{aligned} P(W) &= P(w_1, \dots, w_n) \\ &= P(w_1) \times P(w_2 | w_1) \times \dots \times \\ &\quad P(w_n | w_1, \dots, w_{n-1}) \end{aligned} \quad ,$$

per a una oració donada  $W$ , composta d’una seqüència de tokens de  $w_1$  a  $w_n$ . En termes pràctics, extraïem la NLL mitjana d’una seqüència calculant la seva pèrdua (*loss*; vegeu ‘analysis.py’ al repositori per a detalls d’implementació). Com les parelles d’oracions que comparem només es diferencien al token que les diferencia entre normatives i no normatives, quantifiquem la preferència del model entre oració normativa (a) i oració no-normativa (b) com la diferència entre la NLL mitjana de l’oració (b) i la de l’oració (a). Un valor positiu indica una preferència per la versió normativa, i un valor negatiu, per la no normativa.

Cal esmentar que, mentre que la majoria de parells d’oracions tenen la mateixa longitud, aquelles que impliquen elisió necessàriament difereixen per una paraula (p. ex., caigudes de preposició). Per assegurar-nos que els nostres resultats no estiguin influïts per possibles artefactes

<sup>8</sup><https://huggingface.co/BSC-LT/RoBERTa-ca>

<sup>9</sup><https://huggingface.co/catallama/CataLlama-v0.2-Base>

introduïts per longituds de seqüència diferencials, també vam dur a terme l'anàlisi només amb el subconjunt de frases de longituds iguals. Trobem que els resultats que reportem a §3 s'obtenen independentment d'aquesta qüestió (vegeu l'apèndix per a estimacions numèriques quan només es consideren frases d'igual longitud).

Com que la magnitud de la diferència en NLL depèn del model —entre altres factors, està influïda per la mida del vocabulari— recodifiquem aquesta informació com un resultat binari de preferència a favor o en contra de (a) enfront de (b). Dit d'una altra manera, obtenim 160 judicis per model sobre si assignen una probabilitat més alta a la versió normativa o a la no normativa d'una oració. Abans de continuar a l'avaluació, expliquem breument cadascun dels vuit tipus d'estructures analitzades.

**Estructura 1:** Caiguda de preposició davant la construcció *que*. En registres formals, segons la GIEC (Institut d'Estudis Catalans, 2024, §26.4.1) s'evita el contacte de les preposicions *a*, *de*, *amb*, *en* i la conjunció *que*, de manera que s'elideixen aquestes preposicions.

**Estructura 2:** Preposició *a* davant d'objecte directe. Aquesta construcció és més comuna en castellà, especialment en complements directes animats. Segons la GIEC (Institut d'Estudis Catalans, 2024, §19.3.2), el complement directe en català no va precedit de cap preposició, excepte en casos d'ambigüïtat amb el subjecte o l'alteració de l'ordre bàsic dels constituents oracionals, que no s'han inclòs en aquest corpus d'avaluació.

**Estructura 3:** Preposició *de* introductora de complements regits. Segons la normativa, s'han d'introduir amb la preposició *de*, i no pas *a*, certs complements regits, que per influència del castellà se solen construir malament (Murtra et al., 2025).<sup>10</sup> Alguns exemples normatius són: *contradictori de*, *diferent de*, *fer cas de*, o *rebuig de*.

**Estructura 4:** Alternança de preposicions davant d'infinitiu. Els verbs o altres categories que tenen de complement o adjunt un sintagma nominal introduït per *en* o *amb*, aquestes preposicions alternen amb *a* o *de*, quan introdueixen una subordinada en infinitiu. Tanmateix, segons la GIEC (Institut d'Estudis Catalans, 2024, §26.5.2), la solució preferible en els registres formals és el canvi a *a* o *de*.

**Estructura 5:** Ús no normatiu de *propens de*. Aquesta estructura no és conseqüència directa del contacte amb el castellà, on el verb i la seva preposició són anàlegs a la normativa catalana. La forma normativa en català és *propens a* (Institut d'Estudis Catalans, 2007).

**Estructura 6:** Ús no normatiu de la locució *en motiu de*. Com l'estructura anterior, no és conseqüència del contacte amb el castellà. La forma normativa és *amb motiu de* (Serveis d'Assessorament Lingüístic del Parlament de Catalunya, 2007, Fitxa 7315/2).

**Estructura 7:** Usos no normatius de règims verbals no atribuïbles a la interferència del castellà: *confiar amb*, *continuar amb* i *tractar (de)*. La normativa, estableix *confiar en* (Institut d'Estudis Catalans, 2007), i considera *continuar* un verb transitiu que no admet *amb* (Institut d'Estudis Catalans, 2009). En el sentit de 'prendre alguna cosa com a objecte d'estudi o discussió', la normativa dicta que *tractar* ha de construir-se amb *de* (Institut d'Estudis Catalans, 2009).

**Estructura 8:** Ús no normatiu de la locució *anar en compte*. Com l'estructura anterior, no és conseqüència del contacte amb el castellà. La forma normativa és *anar amb compte* (Institut d'Estudis Catalans, 2024, §19.3.3.2.).

La llista completa d'oracions és disponible en format CSV al nostre repositori. També les reproduïm totes a l'apèndix.

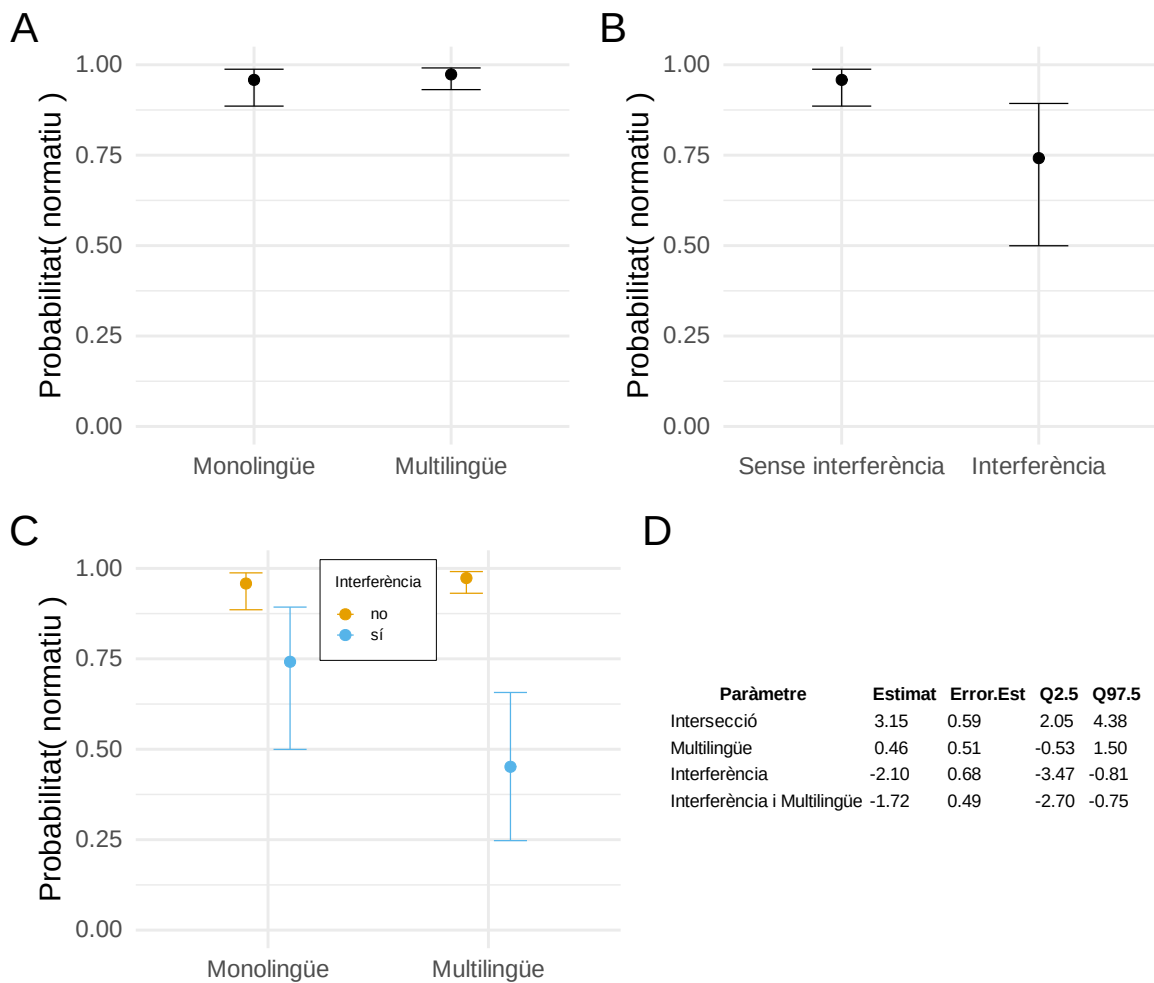
### 3. Resultats

---

A continuació, analitzem la preferència dels models respecte a les nostres parelles mínimes de català (no) normatiu. En particular, quantifiquem si hi ha una diferència entre les preferències dels models monolingües i multilingües; i si les preferències respecte a l'ús normatiu varia en funció de si hi pot haver efectes d'interferència amb el castellà. Per a això, fem una regressió logística amb les 960 preferències obtingudes (160 per model). La regressió estima els efectes principals sobre les preferències per oracions (no) normatives en funció de: (i) la monolingüïtat vs. multilingüïtat de cada model; (ii) possible interferència amb el castellà (oracions del tipus Estructura 1–4 vs. Estructura 5–8); i (iii) la interacció entre aquests dos factors, amb interseccions estimades per oració i per model ("by-sentence" i "by-model" intercepts/efectes aleato-

<sup>10</sup><https://www.uoc.edu/portal/ca/servei-linguistic/criteris/gramatica/preposicions/index.html>





**Figura 1:** Estimacions del nostre model de regressió logística. A: efecte marginal de la monolingüitat vs. multilingüitat sobre la probabilitat de preferir una oració normativa davant la seva alternativa no normativa. B: efecte marginal de preferències entre oracions amb o sense possible interferència amb el castellà. C: interacció entre multilingüitat i interferència. D: estimacions numèriques del model per a cada efecte amb ICs del 95%.

ris).<sup>11</sup> En altres paraules, estimem si el fet que un model sigui multilingüe, sense un ajust final específic per al català, afecta les preferències per les formes normatives; si la interferència amb el castellà les afecta; i, finalment, si hi ha un efecte específic d'interacció d'interferència en el cas dels models multilingües.

El model de regressió bayesià es va estimar fent servir el paquet de R *brms* (Bürkner, 2017) com a interfície amb Stan (Carpenter et al., 2017). Es va inspeccionar per descartar patologies en l'estimació. En particular, vam avaluar la fiabilitat de la mida de la mostra d'estimacions (més de 0,001 mostres efectives per transició) i vam comprovar que el valor de  $R^2$

fragmentat fos inferior a 1,05, que la Bayesian Fraction of Missing Information fos superior a 0,2, i l'absència de trajectòries saturades. La validació creuada presenta valors de pareto- $k < 0,7$ , cosa que suggereix estimacions fiables de la densitat predictiva logarítmica esperada (Vehtari et al., 2017).

La Figura 1 mostra els nostres resultats principals. Els tres primers panells mostren els efectes marginals dels tres factors que estudiem. El quart mostra el resum numèric sencer del model de regressió.

La Figura 1A mostra l'efecte marginal estimat de la multilingüitat sobre la probabilitat de preferir una versió normativa d'una oració. Com s'il·lustra, i també es reflecteix en l'estimació numèrica (Figura 1D), no hi ha una diferència clara entre models mono- i multilingües, en si, pel que fa a la seva preferència per la versió normati-

<sup>11</sup>En sintaxi de *lme4/brms*:  $\text{preferència} \sim 1 + \text{multiling} * \text{interferència} + (1 | \text{oració}) + (1 | \text{model})$  amb "preferència" codificada amb 1 si la diferència entre la NLL no-normativa i la normativa és positiva.

va. En contrast, la Figura 1B mostra l'efecte entre oracions amb interferència (Estructura 1–4) i sense interferència (Estructura 5–8) del castellà. Aquí veiem molta més variació, i el model de regressió suggereix que les oracions amb interferència tenen una probabilitat més baixa de comportar una preferència per la versió normativa. En altres paraules: El castellà interfereix amb el català normatiu. La Figura 1C mostra la interacció marginal entre la multilingüïtat i la interferència: Tot i que a la Figura 1B ja s'observa que els models mostren efectes d'interferència, aquest efecte és encara més pronunciat en els models multilingües. La Figura 1D ofereix un resum numèric de les estimacions de la regressió.

En conjunt, aquests resultats suggereixen que no hi ha una diferència clara entre els models monolingües i multilingües pel que fa a la seva preferència general per les oracions normatives davant alternatives no normatives (Figura 1A i segona fila de la Figura 1D). La interferència amb el castellà disminueix la preferència dels models per l'ús normatiu (Figura 1B i tercera fila de la Figura 1D), i hi ha una disminució addicional de la preferència per la normativitat en els models multilingües en les oracions amb interferència en comparació als monolingües (Figura 1C i quarta fila de la Figura 1D).

#### 4. Discussió

---

Hi ha tres resultats principals en aquesta anàlisi. Primer, a grans trets, els models monolingües i multilingües no difereixen en la seva preferència per oracions normatives. Segon, quan hi ha la possibilitat d'interferència amb el castellà, la preferència de tots dos tipus de models es veu afectada: mostren una major inclinació cap a variants no normatives. Tercer, l'efecte de la interferència difereix entre models monolingües i multilingües: La preferència contra les variants normatives és molt més forta en els models multilingües.

Aquests resultats estan en línia amb les nostres expectatives. Trobem diferències entre models mono- i multilingües allà on són més probables: quan les dades d'altres llengües estretament relacionades, en aquest cas el castellà, apunten en una direcció diferent. No obstant això, és notable que no trobem una diferència per als casos sense interferència. Això vol dir que, per aquests casos i en un domini tancat com les preposicions, els models multilingües desenvolupen competències comparables a les dels models monolingües, independentment que les dades amb què s'han entrenat estiguin desequilibrades (cf. [Brinkmann et al., 2025](#)).

Això ens remet a la discussió de l'apartat 2 relativa a la distinció terminològica entre models monolingües i multilingües. Si més no pel que fa a les dades aquí analitzades, constatem que la distinció té rellevància pràctica, però que la seva rellevància s'ha de jutjar cas per cas. També es podria haver esperat que els models multilingües tinguessin preferències més febles per la norma en els casos sense interferència. Aquest resultat és positiu: els models actualment més populars amb usuaris no mostren desavantatges en almenys alguns aspectes en relació amb models especialitzats en català, que moltes vegades no arriben fins als usuaris.

El segon i tercer resultat —un efecte d'interferència per a tots els models, i particularment per als multilingües— convida a la reflexió i, potser, a l'acció. D'una banda, la variació en dades i entrenament no sembla alterar el fet que la preferència per la normativa disminueixi en contextos d'interferència amb el castellà. Això suggereix que les variants no normatives amb interferència que hem estudiat ja formen part de l'ús del català al qual s'exposen els models. El nostre objectiu aquí no és adjudicar si això hauria de comportar un canvi en la normativa; o bé si caldria reforçar l'ensenyament del català en aquells punts on divergeix del castellà; o bé si convindria preprocessar les dades amb què s'entrenen els models per evitar que l'ús no normatiu sigui propagat pels LLMs. El que sí que mostra clarament aquest resultat és la importància de verificar les preferències lingüístiques dels LLMs, especialment en llengües minoritàries, per tal de tenir aquesta discussió de manera informada; per elaborar polítiques lingüístiques que tinguin en compte l'impacte de les noves tecnologies; i per actuar en conseqüència.

De l'altre, la preferència encara més feble per la norma dels models multilingües demostra com la tecnologia podria actuar com una força de canvi lingüístic: amplificant patrons d'ús i propagant-los. Més en general, aquests resultats mostren que el desenvolupament de tecnologies lingüístiques adaptades a llengües específiques és important i s'hauria de dur a terme conjuntament amb les multilingües. El que funciona per a un conjunt de llengües pot no generalitzar-se a d'altres amb realitats tipològiques i sociolingüístiques diferents. Això és important tant per millorar els models mateixos ([Baucells et al., 2025](#); [Pérez-Mayos et al., 2021](#); [Táboas García et al., 2025](#)) com per detectar conseqüències posteriors potencialment inesperades com les que destaquem aquí.

Una anàlisi complementària que cal fer per posar l'ús lingüístic dels models en context és l'estudi dels parlants de català, cosa que deixem per a futures investigacions. Vam construir aquest corpus d'oracions precisament per estudiar no una diferència entre el que és “correcte” vs. “incorrecte”, sinó entre diferents usos. És ben sabut que molts parlants —catalans o d'altres— no segueixen necessàriament la norma en l'ús que fan de la llengua. Per tant, és probable que fins a cert punt l'efecte d'interferència que detectem és un reflex genuí del canvi lingüístic. Creiem que combinar aquest tipus d'estudi amb experiments conductuals amb subjectes humans és una línia molt prometedora per a futures recerques.

Una altra via prometedora per a futures investigacions és una anàlisi de les dades en brut amb què s'entrenen els models. Aquí ens hem centrat en el comportament lingüístic dels models després del seu entrenament. Però, per deslligar causalment què provoca les seves similituds i diferències, caldria tornar a les dades d'entrada i comparar les taxes de prevalença de les construccions avaluades. Això també aclariria fins a quin punt els LLM amplifiquen o simplement reproduïen els patrons que observen durant l'entrenament; i com interactua aquest fet amb la prevalença de dades d'entrenament mono- vs. multilingües.

## 5. Conclusió

---

Aquest estudi s'ha centrat en l'anàlisi dels biaixos lingüístics dels models de llenguatge en català, avaluant-ne les preferències per construccions normatives enfront de les no normatives, i examinant la influència de la interferència lingüística del castellà. L'avaluació de sis models de llenguatge, tant monolingües com multilingües, a partir d'un corpus de parelles mínimes, revela dos resultats principals.

Primerament, no hem trobat una diferència generalitzada entre els models monolingües i els models multilingües. En segon lloc, la presència de possible interferència amb el castellà redueix la preferència per les variants normatives en tots els models analitzats. Aquesta disminució és marcadament més forta en els models multilingües, la qual cosa podria conduir a una major amplificació d'aquestes preferències en parla catalana generada o modificada per aquests models, que són molt més populars amb usuaris.

Aquests resultats també suggereixen que les variants no normatives influïdes pel castellà estan prou esteses en l'ús real de la llengua com per ser reflectides de manera consistent en les da-

des d'entrenament de tots els models. Donat el paper cada vegada més influent que les tecnologies del llenguatge tenen en la nostra comunicació quotidiana, la comprensió d'aquests biaixos és fonamental. Els models no només reproduïen el llenguatge, sinó que també tenen el potencial d'influir-lo i propagar certs usos. Per tant, aquest estudi subratlla la importància de verificar les preferències lingüístiques dels LLMs per poder fomentar un debat informat.

## Agraïments

---

Agraïm a Iria de Dios Flores, Antoni Oliver González i Lluís Padró pels seus comentaris i suggeriments. TB és finançat pel Ministerio de Ciencia, Innovación, y Universidades, l'Agència Estatal de Investigación i el Fons Social Europeu Plus (RYC2023-045215-I MCIU/AEI/10.13039/501100011033), com també pel projecte EVOSIG PID2024-162668NA-I00 finançat pel MICIU/AEI/10.13039/501100011033/ FEDER, UE.

## Referències

---

- Baucells, Irene, Javier Aula-Blasco, Iria de Dios Flores, Silvia Paniagua Suárez, Naiara Perez, Anna Salles, Susana Sotelo Docio, Júlia Falcão, Jose Javier Saiz, Robiert Sepulveda Torres, Jeremy Barnes, Pablo Gamallo, Aitor Gonzalez-Agirre, German Rigau & Marta Villegas. 2025. IberoBench: A benchmark for LLM evaluation in Iberian languages. En *31<sup>st</sup> International Conference on Computational Linguistics (COLING)*, 10491–10519. [↗](#)
- Brinkmann, Jannik, Chris Wendler, Christian Bartelt & Aaron Mueller. 2025. Large language models share representations of latent grammatical concepts across typologically diverse languages. En *Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics (NAACL)*, 6131–6150. [doi](#) [10.18653/v1/2025.naacl-long.312](https://doi.org/10.18653/v1/2025.naacl-long.312)
- Bürkner, Paul-Christian. 2017. brms: An R Package for Bayesian Multilevel Models using Stan. *Journal of Statistical Software* 80(1). [doi](#) [10.18637/jss.v080.i01](https://doi.org/10.18637/jss.v080.i01)
- Carpenter, Bob, Andrew Gelman, Matthew D. Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li & Allen Riddell. 2017. Stan: A probabilistic programming language. *Journal of Statistical Software* 76(1). [doi](#) [10.18637/jss.v076.i01](https://doi.org/10.18637/jss.v076.i01)

- Gonzalez-Agirre, Aitor, Marc Pàmies, Joan Llop, Irene Baucells, Severino Da Dalt, Daniel Tamarayo, José Javier Saiz, Ferran Espuña, Jaume Prats, Javier Aula-Blasco, Mario Mina, Iñigo Pikabea, Adrián Rubio, Alexander Shvets, Anna Sallés, Iñaki Lacunza, Jorge Palomar, Júlia Falcão, Lucía Tormo, Luis Vasquez-Reina, Montserrat Marimon, Oriol Pareras, Valle Ruiz-Fernández & Marta Villegas. 2025. Salamandra technical report. arXiv [cs.CL]. [doi 10.48550/arXiv.2502.08489](https://doi.org/10.48550/arXiv.2502.08489)
- Institut d'Estudis Catalans. 2007. Diccionari de la llengua catalana [diec2, online version]. [↗](#)
- Institut d'Estudis Catalans. 2009. Manual d'estil. la redacció i l'edició de textos (4a ed.). [↗](#)
- Institut d'Estudis Catalans. 2024. *Gramàtica de la llengua catalana [GIEC, online version]*. Barcelona: IEC. [doi 10.2436/10.2500.08.1](https://doi.org/10.2436/10.2500.08.1)
- Le Scao, Teven, Angela Fan, Christopher Akiki, Ellie Pavlick & et al. 2023. BLOOM: A 176b-parameter open-access multilingual language model. arXiv [cs.CL]. [doi 10.48550/arXiv.2211.05100](https://doi.org/10.48550/arXiv.2211.05100)
- Liang, Weixin, Yaohui Zhang, Zhengxuan Wu, Haley Lepp, Wenlong Ji, Xuandong Zhao, Hancheng Cao, Sheng Liu, Siyu He, Zhi Huang, Diyi Yang, Christopher Potts, Christopher D. Manning & James Zou. 2025. Quantifying large language model usage in scientific papers. *Nature Human Behaviour* [doi 10.1038/s41562-025-02273-8](https://doi.org/10.1038/s41562-025-02273-8)
- Lin, Xi Victoria, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, Ramakanth Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav Chaudhary, Brian O'Horo, Jeff Wang, Luke Zettlemoyer, Zornitza Kozareva, Mona Diab, Veselin Stoyanov & Xian Li. 2022. Few-shot learning with multilingual language models. arXiv [cs.CL/cs.AI]. [doi 10.48550/arXiv.2112.10668](https://doi.org/10.48550/arXiv.2112.10668)
- Murtra, Pilar, Thomas Bell, David Cullen, Jordi Gavaldà, Cristina López, Xavier Marzal & Alba Pérez. 2025. Servei lingüístic de la Universitat Oberta de Catalunya. [↗](#)
- Papadimitriou, Isabel, Kezia Lopez & Dan Jurafsky. 2023. Multilingual BERT has an accent: Evaluating English influences on fluency in multilingual models. En *Findings of the Association for Computational Linguistics*, 1194–1200. [doi 10.18653/v1/2023.findings-eacl.89](https://doi.org/10.18653/v1/2023.findings-eacl.89)
- Pérez-Mayos, Laura, Alba Táboas García, Simon Mille & Leo Wanner. 2021. Assessing the syntactic capabilities of transformer-based multilingual language models. En *Findings of the Association for Computational Linguistics*, 3799–3812. [doi 10.18653/v1/2021.findings-acl.333](https://doi.org/10.18653/v1/2021.findings-acl.333)
- Schut, Lisa, Yarin Gal & Sebastian Farquhar. 2025. Do multilingual LLMs think in English? En *ICLR 2025 Workshop on Building Trust in Language Models and Applications*, [↗](#)
- Serveis d'Assessorament Lingüístic del Parlament de Catalunya. 2007. Optimot, consultes lingüístiques. [↗](#)
- Táboas García, Alba, Piotr Przybyła & Leo Wanner. 2025. Exploring morphology-aware tokenization: A case study on Spanish language modeling. En *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 30493–30506. [doi 10.18653/v1/2025.emnlp-main.1552](https://doi.org/10.18653/v1/2025.emnlp-main.1552)
- Vehtari, Aki, Andrew Gelman & Jonah Gabry. 2017. Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Statistics and Computing* 27. 1413–1432. [doi 10.1007/s11222-016-9696-4](https://doi.org/10.1007/s11222-016-9696-4)
- Wendler, Chris, Veniamin Veselovsky, Giovanni Monea & Robert West. 2024. Do Llamas work in English? on the latent language of multilingual transformers. En *62<sup>nd</sup> Annual Meeting of the Association for Computational Linguistics*, 15366–15394. [doi 10.18653/v1/2024.acl-long.820](https://doi.org/10.18653/v1/2024.acl-long.820)

## Apèndix

**Comprovacions de robustesa** La Taula 1 mostra les estimacions d'un model de regressió que inclou CataLlama-v0.2 i RoBERTa-ca. L'estructura del model de regressió i tots els altres detalls són els mateixos que en el text principal. Per a RoBERTa-ca, la (pseudo-)negative log likelihood es calcula prenent el logaritme del likelihood de les oracions amb les preposició clau emascarades.

Com es pot observar, les estimacions són semblants a les obtingudes sense aquests dos models (comparar amb la Figura 1D). Principalment, l'efecte d'interferència és més gran amb la inclusió d'aquests dos models. Això suggereix que els resultats són relativament robustos davant de la variació en el conjunt de models estudiats, almenys pel que fa als monolingües.



Paràmetre	Estimat	Error.Est	Q2.5	Q97.5
Intersecció	3.45	0.51	2.51	4.56
Multilingüe	0.16	0.42	-0.64	0.96
Interferència	-2.72	0.60	-3.99	-1.62
Interferència i Multilingüe	-0.99	0.40	-1.82	-0.22

**Taula 1:** Estimacions numèriques per a un model de regressió que inclou CataLlama i RoBERTa-ca.

Paràmetre	Estimat	Error.Est	Q2.5	Q97.5
Intersecció	3.05	0.55	2.03	4.16
Multilingüe	0.69	0.56	-0.38	1.84
Interferència	-1.40	0.57	-2.55	-0.33
Interferència i Multilingüe	-1.59	0.52	-2.64	-0.59

**Taula 2:** Estimacions numèriques per a un model de regressió amb només estímuls amb el mateix nombre de paraules per parella (no) normativa.

La Taula 2 mostra les estimacions d'un model de regressió quan només es fan servir oracions d'igual longitud. S'obtenen les mateixes tendències que a la Fig. 1, cosa que indica que les longituds de seqüència diferencials ocasionals (p. ex., caigudes de preposició) no afecten els nostres resultats de maneres inesperades.

**Estímuls** A continuació es mostren els parells de oracions que vam analitzar. Com en els exemples del text principal, l'opció normativa és la primera dins dels claudàtors, abans de la coma. Les entrades buides indiquen que la respectiva variant (no) normativa deixa aquella posició buida.

### Estructura 1

1. No m'havia adonat {, de} que estava mal col·locat.
2. Abans {, de} que comencés, ja ho sabíem.
3. Depèn {, de} que decideixin els experts.
4. No es recorda {, de} que li vaig explicar.
5. No m'havia adonat {, de} que estava equivocat.
6. Depèn {, de} que vulgin fer.
7. No es recorda {, de} que va veure aquell dia.
8. Tinc la certesa {, de} que arribaràs a temps.
9. M'alegro {, de} que ho sàpigues.
10. Estic convençut {, de} que tenim raó.
11. Aspirava {, a} que el triessin.
12. Contribuirà {, a} que el projecte tiri endavant.
13. S'oposen {, a} que la normativa es modifiqui.
14. Em refereixo {, a} que va fer una nova pel·lícula.

15. Em refereixo {, a} que va cantar.
16. Estic acostumat {, a} que m'avisin amb temps.
17. L'autoritat obliga {, a} que es pagui un peatge.
18. N'hi ha prou {, amb} que diguis la veritat.
19. Confio {, en} que tot sortirà bé.
20. Insisteix {, en} que tenim una altra opció.

### Estructura 2

1. Això afecta {el, al} pintor solidari.
2. He vist {el, al} senyor calb.
3. L'entrenador felicitarà {els, als} jugadors després del partit.
4. L'empresa vol formar {els, als} seus treballadors.
5. Estic esperant {el, al} professor per resoldre un dubte.
6. Han acomiadat {el, al} empleats més antics de l'empresa.
7. Sempre ajuda {els, als} seus companys quan ho necessiten.
8. El professor ha avisat {els, als} alumnes sobre el canvi d'horari.
9. El jutge ha interrogat {els, als} testimonis del cas.
10. Han seleccionat {els, als} millors candidats per al lloc de feina.
11. La policia busca {el, al} culpable.
12. El jutge va absoldre {els, als} acusats.
13. Necessito trobar {el, al} metge.
14. El soroll va despertar {el, al} veí.
15. Vam citar {el, al} director de l'escola.
16. La notícia ha afectat {el, al} jove empresari.
17. Necessiten substituir {el, al} tècnic que ha marxat.
18. El llibre va impressionar {els, als} lectors més joves.
19. L'entrenador ha sancionat {el, al} jugador titular.
20. Van ajudar {els, als} excursionistes perduts.

### Estructura 3

1. Al restaurant fa molta olor {de, a} fregit.
2. Tinc por {de, a} les aranyes.
3. El meu horari és diferent {del, al} teu.
4. La solució {del, al} problema és molt senzilla.
5. La noia va aprofitar per treure profit {de, a} la situació i va guanyar molts diners.
6. Va fer cas {de, a} les recomanacions del metge.
7. Sento una forta pudor {de, a} cremat a la cuina.
8. L'article fa esment {de, a} diverses investigacions recents.
9. Va mostrar un clar rebuig {de, a} la proposta.
10. Hauries de fer un bon repàs {de, a} la matèria abans de l'examen.

11. La seva versió és contradictòria {de, a} la que vam escoltar.
12. La pel·lícula és diferent {de, a} la novel·la.
13. No facis cas {de, a} les males notícies.
14. El presentador va fer esment {de, a} la polèmica.
15. Sentia una forta olor {de, a} tabac apagat.
16. Els nens tenien por {dels, als} gossos grans.
17. La noia porta un vestit {de, a} ratlles blanques i negres.
18. El jurat fa una menció especial {de, a} l'obra d'aquesta estudiant.
19. Això no implica cap increment {de, a} les {de, a}speses {de, a}l web.
20. Sempre ha volgut estudiar una carrera diferent {de, a} la que ha fet.

#### Estructura 4

1. Tenim molt d'interès {a, en} complir els nostres compromisos.
2. La idea consisteix {a, en} analitzar els models.
3. La majoria d'estudiants es van mostrar d'acord {a, en} fer una pausa.
4. Comptava {a, amb} fer això.
5. Insisteixen {a, en} tornar cap a casa.
6. Va mostrar interès {a, en} aprendre nous mètodes per fer aquests exercicis.
7. L'exercici es basa {a, en} explicar les circumstàncies del conflicte.
8. Ens complaem {a, en} acceptar la vostra col·laboració.
9. Engrescat {a, en} fundar la coral, tot el temps que hi dedica li sembla poc.
10. Estava tan capficat {a, en} completar la jugada, que no es va adonar que el temps s'havia acabat.
11. Tinc interès {a, en} aprendre més llengües.
12. El treball consisteix {a, en} establir una relació entre les dues variables.
13. Compte {a, amb} demanar el mateix a tots dos llocs
14. L'empresa es va comprometre {a, amb} lliurar la mercaderia abans de la data.
15. La subsistència es basava {a, en} intercanviar productes.
16. S'ha entossudit {a, en} comprar folis reciclats.
17. La democràcia consisteix {a, en} fer que el poble participi en la política.
18. No s'ha de capficar {a, en} recordar totes les dades de memòria.
19. Cal invertir un gran esforç {a, en} millorar l'atenció al client.
20. Ell es complau {a, en} resoldre trencaclosques complicats.

#### Estructura 5

1. És molt propens {a, de} agafar refredats a l'hivern.
2. Aquesta planta és propensa {a, de} les plagues si no la cuides bé.
3. Els adolescents són més propensos {a, de} prendre riscs innecessaris.
4. Les persones amb pells clares són propenses {a, de} les cremades solars.
5. El meu oncle és propens {a, de} explicar sempre les mateixes anècdotes.
6. Aquesta zona és propensa {a, de} les inundacions quan plou fort.
7. Som propensos {a, de} oblidar les claus a casa.
8. No sóc gens propens {a, de} enfadar-me per tonteries.
9. Les economies inestables són propenses {a, de} patir crisis financeres.
10. Aquest material és propens {a, de} l'oxidació si s'exposa a la humitat.
11. El meu veí és molt propens {a, de} exagerar les històries.
12. La pell sensible és propensa {a, de} reaccions al·lèrgiques.
13. Aquella noia és propensa {a, de} dir mentides sense motiu.
14. Els cotxes vells són propensos {a, de} avaries inesperades.
15. És propens {a, de} cometre errors quan està sota pressió.
16. Ets massa propens {a, de} canviar d'opinió a l'últim moment.
17. Aquesta màquina és propensa {a, de} errors si no es calibra bé.
18. El meu soci és propens {a, de} prendre decisions arriscades.
19. La nostra regió és propensa {a, de} grans nevades a l'hivern.
20. Les persones nervioses són propenses {a, de} mossegar-se les ungles.

#### Estructura 6

1. {Amb, En} motiu de la seva jubilació, li van organitzar una festa sorpresa.
2. La botiga va tancar {amb, en} motiu de les obres al carrer.
3. Es va organitzar un concert benèfic {amb, en} motiu de la recaptació de fons.
4. {Amb, En} motiu del Dia del Llibre, les llibreries fan descomptes.
5. Van cancel·lar el vol {amb, en} motiu del mal temps.
6. L'ajuntament va emetre un comunicat {amb, en} motiu de la nova normativa.

7. Es va tallar el trànsit {amb, en} motiu de la manifestació.
8. La reunió es va ajornar {amb, en} motiu de l'absència del director.
9. {Amb, En} motiu de la celebració, el museu ofereix entrada gratuïta.
10. L'empresa va fer una donació {amb, en} motiu de la campanya solidària.
11. {Amb, En} motiu de les festes, l'ajuntament va tancar.
12. La trobada es va convocar {amb, en} motiu de la crisi energètica.
13. Es va avançar l'hora {amb, en} motiu de la celebració del partit.
14. L'excursió es va cancel·lar {amb, en} motiu de la forta ventada.
15. Vam assistir a la conferència {amb, en} motiu de la presentació del llibre.
16. Han declarat dia festiu {amb, en} motiu del patronatge del poble
17. Vam fer una excepció {amb, en} motiu de la seva situació particular.
18. Vam llançar una edició especial {amb, en} motiu del centenari de la marca.
19. La subvenció va ser concedida {amb, en} motiu de les bones previsions.
20. {Amb, En} motiu dels preparatius, la pista d'esquí roman tancada.
14. Continuem {, amb} la tasca que vam deixar a mig fer.
15. El debat tractarà {de, } les noves propostes de llei.
16. En aquest llibre l'autor tracta molt sintèticament {de, } la morfologia.
17. El més interessant del llibre és com tracta {de, } la sintaxi.
18. L'informe no tracta {de, } les qüestions més importants.
19. La pel·lícula tracta {de, } la vida de Mozart.
20. La conferència tractava {de, } la intel·ligència artificial.

### Estructura 8

### Estructura 7

1. Pots confiar {en, amb} mi per guardar el secret.
2. No confio {en, amb} les seves promeses, sempre les incompleix.
3. Cal confiar {en, amb} el criteri dels experts.
4. He après a no confiar {en, amb} qualsevol que sembli amable.
5. Perquè la relació funcioni, heu de confiar l'un {en, amb} l'altre.
6. Confiem {en, amb} que tot sortirà bé al final.
7. És difícil confiar {en, amb} algú que t'ha mentit abans.
8. Els ciutadans confien {en, amb} la justícia del seu país.
9. Si us plau, continueu {, amb} la lectura del proper capítol.
10. El govern continuarà {, amb} les negociacions la setmana vinent.
11. Després d'aquesta interrupció continuem {, amb} el programa.
12. Malgrat tot, l'equip continua {, amb} aquest projecte.
13. Podem continuar {, amb} la feina demà al matí.
1. Vés {amb, en} compte quan creuis el carrer.
2. Aneu {amb, en} compte amb el terra, que rellisca.
3. Cal anar {amb, en} compte a l'hora de signar contractes.
4. Sempre vaig {amb, en} compte amb el que dic en públic.
5. Van anar {amb, en} molt de compte per no fer soroll.
6. Si no vas {amb, en} compte, et pots fer mal.
7. Anem {amb, en} compte, que no sabem què ens trobarem.
8. El metge li va dir que anés {amb, en} compte amb els greixos.
9. És una situació delicada, hem d'anar {amb, en} compte.
10. Aneu {amb, en} compte amb les estafes per internet.
11. El metge li va dir que anés {amb, en} compte amb la dieta.
12. Sempre va {amb, en} compte a l'hora d'invertir els seus diners.
13. L'advocat els va aconsellar anar {amb, en} compte amb les clàusules del contracte.
14. Ves {amb, en} compte quan entris al bosc, podria haver-hi animals.
15. El cuiner va haver d'anar {amb, en} compte amb l'oli calent.
16. Ves {amb, en} compte amb qui comparteixes la informació personal.
17. Si vas {amb, en} compte, evitaràs la majoria dels problemes.
18. Ella sempre va {amb, en} compte amb el que escriu a les xarxes socials.
19. El doctor va advertir-li que anés {amb, en} compte amb l'estrès.
20. Cal anar {amb, en} compte en les declaracions públiques.