# French partial interrogatives: a microdiachronic corpus study of variation and new perspectives in a refined pragmatics framework

**Gabriel Thiberge**[*], **Flora Badin**[**], **Loïc Liégeois**[***]

## Introduction

This paper presents new elements for the understanding of the variation exhibited in French with regards to the syntactic realization of partial interrogatives. The various interrogative variants available to French native speakers (e.g., 1a-f[1]) have often been analysed as structurally derived from the canonical S(ubject)-V(erb)-O(bject) structure of declarative sentences (2a) in French where OVS order is ungrammatical (2b). Under this view, there is a Derivational Complexity Metrics (DCM; see Jakubowicz, 2011) where a structure like (1a) is less complex to produce or understand than, say, (1b) with at the same time fronting of the interrogative element (Wh- phrase) and inversion of subject-verb order. Very broadly speaking there is a crucial opposition between the so-called *in situ* variant (1a) and other forms where the Wh- element is not in this position (hence sometimes called *ex situ*). In these configurations, further complexity is added with each transformation of the "base" structure, be it the "est-ce que" insertion (Fesk) or a verb-subject inversion (FINV) for instance. (1e) and (1f) are other interrogative variants where a binary opposition between *in-* and *ex- situ* positions is less clear (see Hamlaoui (2008, 2009) or Clech-Darbon et al. (1999) for more in-depth analyses of the cleft structure).

[**] LLF, UMR 7110, Université de Paris, Sorbonne Paris Cité, CNRS, gtac@tuta.io

[***] LLL, UMR 7270, Université d'Orléans, COMUE Centre-Val de Loire, CNRS,

[****] LLF, UMR 7110, Université de Paris, Sorbonne Paris Cité, CNRS et CLILLAC-ARP, EA 3967, Université de Paris, Sorbonne Paris

[1]     Those examples are in no way an exhaustive list of all variants, see among others Coveney (2011) or Delaveau (to appear) for more complete presentations of partial interrogation and, more broadly, of interrogatives in French.

(1)   a.    Tu        vois            qui?[2]
                 2SG      see.PST.2SG    who

         b.    Qui    vois-tu?[3]
                 who    see.PST.2SG-2SG

         c.    Qui    tu        vois?[4]
                 who    2SG        see.PST.2SG

         d.    Qui    est-ce que    tu           vois?[5]
                 who    INTEXP       2SG         see.PST.2SG

         e.    C'est           qui    que    tu    vois?[6]
                 EXPL.be.PST.3SG   who    COMP    2SG    see.PST.2SG

         f.    Qui    que        tu           vois?[7]
                 who    COMP       2SG         see.PST.2SG

  Who do you see?

(2)   a.    Tu       vois         quelqu'un.
                 2SG      see.PST.2SG    someone

         b.    *Quelqu'un    tu        vois.
                 someone     2SG       see.PST.2SG

This very concept of gradual complexity was the basis for a number of acquisition studies on interrogative sentences in French. Those works tried to put forth an explanation as to why native French speaking children seem to produce more *in situ* sentences at the beginning of their linguistic development, with variants like (1b) or (1c) generally appearing much later (Hulk, 1996). These studies based on the idea of syntactic complexity primarily focused on French (among others: Zuckermann & Hulk, 2001; Hamann, 2006; Strik, 2006; Jakubowicz & Strik, 2008[8]), but their scope also extended to other languages where multiple interrogative variants are available, such as Portuguese (among others: Soares, 2003, 2006; Baião & Lobo, 2014).

Other linguistic levels of analysis were also explored on this alternation phenomenon. Hamlaoui (2009, 2010), for example, found corpus-based evidence for an influence of phonotactics on the choices made by French native

---

[2]    In situ (IS)

[3]    Fronting+Inversion (FINV)

[4]    Simple Fronting (F)

[5]    Fronting+ESK (Fesk)

[6]    Cleft

[7]    Fronting + Complementizer

[8]    See Prévost (2009) for a broader compilation of data on this question.

speakers in their production of partial interrogatives. She established a correlation between the length (in syllables) of the non-interrogative part of the sentence and the likelihood of fronting of the Wh-phrase. In a more pragmatics-oriented line of work, Boeckx (1999), Beyssade (2006) or Déprez & al. (2013), among others, linked *in situ* interrogatives to focus positions, which Hamlaoui (2009, 2010, 2011) disputes. In line with other works (Coveney, 1995, 1997; Krifka, 2007 ; Engdhal, 2006; a.o.), Hamlaoui argues for an influence of the non-Wh part of the interrogative sentence, depending on whether it is given in the immediate discourse or situation. Adding yet another perspective, sociolinguistic works from as early as the 1960's established a bridge between sociolinguistic groups and the favored use of interrogative variants. For instance, Ashby (1977) concludes that academics use more FINV sentences, with both fronting of the Wh-element and verb-subject inversion, than other social groups he studied (see Quillard, 2001, for more examples of the same nature).

In this paper, we present new data allowing for a pragmatically refined sociolinguistic approach to the variation phenomenon. We will discuss our data in relation to "3rd-wave" sociolinguistics (Eckert, 2012), probabilistic pragmatics (Goodman & Lassiter, 2015) and, by extension, to Lewis' Game Theory (Lewis, 1969). Variation is no longer just used to mark social or cultural belonging, and rather becomes a tool to socially position oneself and to adapt this position during the linguistic exchange. This framework is inspired by the latest developments on the social meaning of utterances (Burnett, 2017): Speakers behave differently according to their goals and to the persona (Ochs, 1992) or "social mask" they want to convey during interaction, and the nature of interaction impacts strategies (i.e. context formality/colloquiality or presence/absence of an audience to the exchange can have consequences on what is said and how it is said).

The data we will present are from the two subparts of the ESLO corpus of oral French (Eshkol-Taravella & al., 2012; LLL, 2017): ESLO1, recorded in the 1960s, and ESLO2, recorded with the same protocols in the 2010's. The first crucial step in our study was to prepare the data for an import in the TXM tool (Heiden & al., 2010) and a transposition of sociolinguistic metadata for all recorded speakers on the interrogative sentences we extracted (SECTION 1). The two periods of recording allow for a study of diachronic evolution of French (i.e. French speakers from the 1960s did not use the same variants as speakers from the 2010's, in similar contexts and for similar social groups, SECTION 2). In the same section we

then show that the great diversity of speakers in the ESLO project also allows for an exploration of age factors, with speakers aged 15-25 behaving differently from 35-55 year-old speakers, across interaction contexts and time periods. In SECTION 3, the high variability of recording contexts in the ESLO project allows us to show the importance of context factors: for the same age group and over the same time window, French speakers behave differently in different social contexts.

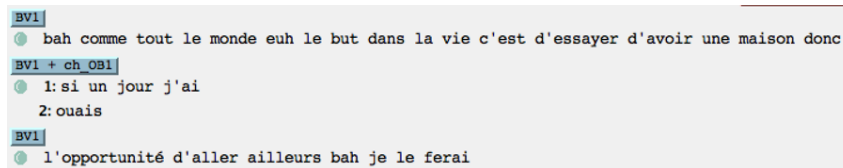## 1.  Data, corpus extraction and sociolinguistic metadata

The ESLO corpus (Eshkol-Taravella et al., 2012 ; LLL, 2017) is a database composed of transcripts, recordings and metadata. This spoken corpus is divided into two subparts: the first one, ESLO1, is based on data recorded in the 1960s while the second one, ESLO2, is based on data recorded with the same protocol since 2010 (data collection still in progress). Based, respectively, on 3.1 million tokens (for 280 hours of spoken data recorded) for the first part and 1.8 million tokens (for 142 hours of spoken data recorded) for the second part, the two subparts of the ESLO corpus are available for research purposes at no cost. The main value of this corpus is to have associated rich metadata, both at recording and speaker levels. For example, each recording is associated with information about the interaction situation (interview, conference, natural interactions…). As for speakers, access is given to sociolinguistic characteristics like age, education or socio-professional group.

Even though the ESLO team proposes an online concordancer, it does not exploit all these metadata and does not allow requests using part of speech tagging. In order to analyse sociolinguistic variation concerning various syntactic structures of interrogatives, we chose to exploit the ESLO corpus with a reference textometric tool, TXM (Heiden et al., 2010). Among other qualities, this tool executes, during the data importation process, a lemmatisation and a part-of-speech tagging, based on a Tree-Tagger (Schmid 1994, 1995). Moreover, TXM allows to create several sub-corpora based on metadata features. Before using this tool and exploiting all of these specificities, we had to prepare the data for importation.

After downloading the Transcriber (Barras et al., 2001) version of the corpus, we set up a series of processes to convert the native XML format of the corpus to an XML format that could be imported into TXM. An automatic conversion from Transcriber XML format to TXM is possible

either with the "Importation" function of TXM or with the TEI-CORPO tool (Liégeois et al., 2015 ; MoDyCo, 2016), but we had to develop our own script in order to manage two important settings of our study: the conservation of the transcribed segments' continuity and the projection of metadata at utterance level.

Concerning the transcribed segments continuity, a problem came from ESLO's transcription norms. Indeed, the transcripts are segmented neither at the utterance level (as defined by Parisse & Le Normand, 2006, for example) nor at speech turns but according to discursive criteria like speech overlapping. Thus, an utterance can be segmented on several tiers (Fig. 1), which can affect the quality of part-of-speech tagging. In the example below, the production of the speaker "BV1" is segmented in three tiers, although it corresponds to a unique speech turn and a unique utterance, as defined in Parisse & Le Normand (2006). We solved this problem by automatic reconstruction of speech turns impacted by overlapping. After this, the production of the speaker "BV1" below is structured on a unique tier that corresponds to one speech turn.r an example.



*figure 1*:  Extract from an overlap in the transcript ESLO2_ENT_1001

Regarding metadata projection, we were able to rely on the possibilities offered by TXM for data importation. We created a spreadsheet file containing for each recording (represented in rows), a set of speaker metadata (represented in columns). During the importation process, the tool automatically inserted these features in XML arguments at the speech turn level, thus each speech turn is documented with speaker information concerning their unique ID, sex, age, profession, socio-professional category, education.

All these metadata could then be used to conduct contrastive studies by creating sub-corpora on the basis of one or more variables. To these was added information relating to the interaction situation (interviews, natural interaction during meal-time, scientific conference...) and information on the period of data collection (ESLO1 or ESLO2). It is also at the time of importation that part-of-speech tagging was done, using TreeTagger and a model especially developed for French oral data by the PERCEO

project (Benzitoun et al., 2012 ; ATILF, 2012). In the end, we obtained a version of the ESLO corpus that could be used for interrogative extraction using discursive and socio-linguistic criteria.

To extract the data which were useful for our study, we used the TXM query engine by performing a CQL (Corpus Query Language) query to obtain all the interrogatives of the targeted sub-corpora. Regarding the data extracted from interviews, we excluded all of the interviewers' productions. Indeed, interviewers are non-naive speakers, in the sense that they are linguists and/or French-language teachers. But, above all, interviewers of the ESLO1 sub-corpus are not native speakers of French since this part of the corpus was collected as part of a research project carried out by British academics.

## 2. A microdiachronic view of French Partial Interrogatives (FPIs), from the 1960s to the 2010s

### 2.1. Extraction and annotation

Automatic extraction by the methodology described above returned a high number of tokens. Among those, a manual annotation for syntactic features was conducted by two different annotators. The goal here was to manually check the structure of FPIs, and to exclude relative clauses beginning by "qui" and "que" (homophonic to Wh- words) and other non-relevant tokens. The syntactic criteria used were of binary nature: either they were met by the token(=1), or they were not met (=0). The criteria were:
– FPI: is the token asking for a missing piece of information about the world (1) or not (0) (i.e. excludes relative clauses but also speech acts like demands for repetition, etc.);
– Solo: is there a verb directly overseeing the Wh-element (0) or not (1);
– Root: is the token a root clause (1) or not (0);
– Embedded: is the Wh- element in the matrix part of the sentence (0) or not (1);
– Infinitive: is the Wh- element an argument of a finite verb (0) or not (1);
– Fronted: is the Wh- element at the beginning of the sentence (1) or not (0);
– Inversion: is there a verb-subject inversion (1) or not (0) in the verb phrase the Wh-element is an argument of;
– Esk: is the idiomatic interrogative expression "est-ce que" connecting the Wh-element to the rest of the interrogative clause (1) or not (0);

– NEGATION: is the verb directly overseeing the Wh-element overtly negated (1) or not (0); double negation (ne... V... pas) and simple negation (ø ... V... pas) were coded in the same way.

A list of the Wh- elements used in the interrogatives was also compiled.

Subject interrogatives, where the Wh- element is the syntactic subject of the verb, were also annotated as such, because they constrain word order (no *ex situ* configuration possible).

All in all, for both the 15-25 yo and 35-55 yo age categories in both publicly available corpora (ESLO1 & ESLO2), 1715 tokens were automatically extracted in TXM. It then became obvious that some investigators from the ESLO2 project were coded in the metadata in the same way as linguistic consultants. The tokens from those non-naive speakers were excluded and only 1399 from naive speakers were kept for annotation. Among those 1399 tokens, 617 root FPIs forming a complete finite sentence were kept for statistical analysis, from 130 different transcription files, each file corresponding to a different recording at the time of data collection.

Criteria for exclusion were:
– not a partial interrogative (= for example, extractions where "qui" was a relativizer);
– not a full sentence (= for example, a wh- word being used in its own);
– not a root interrogative (=embedded clause, constrains the word order);
– not a finite interrogative (high probablity of not having an overt subject in the sentence);
– subject interrogatives (constrains the word order).

## 2.2.    Overall distribution of French Partial Interrogatives and statistical tools

The 617 FPIs kept for analysis are distributed in different proportions over different age and corpus subsets relevant to our analysis, as summarized in Table 1.

*table 1*:    Distribution of FPIs by age group and corpus

|  | **ESLO1** | **ESLO2** | **TOTAL** |
|---|---|---|---|
| 15-25 y.o. | 60 | 88 | 148 |
| 35-55 y.o. | 192 | 277 | 469 |
| Total | 252 | 365 | 617 |

The figures in Table 1 show the collected data are far from perfectly balanced across subsets, and most notably there seems to be a lack of FPIs collected from the 15-25 groups, as compared to 35-55 groups, regardless of the time window. The statistical modelling tools used for the analyses presented in the next sections take those numeric disparities into account.

The FRONTED, ESK, and INVERSION syntactic criteria were used to compute the proportions of use for the four main interrogative types of sentences given in (1a-d) and reproduced below.

(1)     a.     Tu        vois              qui?[9]
               2SG       see.PST.2SG       who
        b.     Qui      vois-tu?[10]
               who      see.PST.2SG-2SG
        c.     Qui      tu          vois?[11]
               who      2SG         see.PST.2SG
        d.     Qui      est-ce que    tu            vois?[12]
               who      INTEXP        2SG           see.PST.2SG

Table 2 gives an overview of the findings in raw numbers, on which we based our analyses. The first notable fact here is the obvious prevalence of simple fronting and *in situ* forms overall, as compared to two other kinds of fronting (either with the "est-ce que" idiomatic expression or with verb-subject inversion).

*table 2*:   Distribution of FPIs by age group, corpus and position of the Wh- element

|            | ESLO1 |      |      |     | ESLO2 |      |      |     |
|------------|-------|------|------|-----|-------|------|------|-----|
|            | f     | fesk | finv | is  | f     | fesk | finv | is  |
| 15-25 y.o. | 24    | 13   | 8    | 15  | 16    | 6    | 6    | 60  |
| 35-55 y.o. | 68    | 30   | 47   | 47  | 74    | 28   | 24   | 151 |
| Total      | 92    | 43   | 55   | 55  | 90    | 34   | 34   | 211 |

The next sections make use of two very different types of statistical analyses: simple descriptive frequencies on the one hand to provide a very broad idea of the observable phenomena, and Bayesian regression modelling on the other hand, which allows for

---

[9]      In situ (IS)

[10]     Fronting+Inversion (FINV)

[11]     Simple Fronting (F)

[12]     Fronting+ESK (Fesk)

generalizations based on the raw observed frequencies by measuring the effect strength of some predictors (age group, time period, context of interaction) in the observed data. The choice of a Bayesian rather than other more classical frequentist approaches relies on several factors.

First, Bayesian modelling allows for a non-binary take on the data, without relying on the opposition between significant and non-significant effects. Rather, it is a fine-grained analysis of the distribution of the probability of the effects being real, in relation to the potential priors (hypothesis) set on these distributions and to the new data fed to the model (see Sorensen & al. 2016, for a more detailed account on why Bayesian statistics can be a useful tool for linguistic data exploration). The crucial point here is this kind of modelling can account for data presenting a high degree of variance, derived from small data sets, which is of particular interest here.

A second advantage of the "Bayesian way" is of a more practical nature. The data manipulated here is neither continuous (like time-related data) nor binomial (binary variable with only two levels), but mostly categorical, both for the dependent variable (DV, here the type of interrogatives) and the predictors tested (age group, time period, context). As such, the most suitable models available would have been vector general linear models (VGLM's, Yee & Wild, 1996; Yee, 2015), but those models cannot easily include random effects, such as the potential variability between all the transcriptions in the corpora (as a reminder:130 files for 617 tokens in total). Bayesian modelling with R packages like stan-based brms (R Development Core Team, 2009; Carpenter & al., 2017; Bürkner, 2017; Bürkner, 2018), allows this kind of refining, and was thus selected.

## 2.3.   A diachronic change in the use of interrogative structures

From a first superficial perspective, basic frequencies computed from Table 2 shed light on an overall change in the interrogative habits between the two time periods (Fig. 2). In the ESLO1 corpus, covering the 1960's, more than 75% of the extracted questions have a fronted Wh-element, with nearly a half being simple frontings (F). In the ESLO2 corpus, covering the 2010's, nearly 60% of the extracted tokens have the Wh-element *in situ* (IS), with drastic reductions from all three kinds of fronting. Frontings with an additional element (Fesk or FINV) seem to have reduced more drastically (by slightly more or slightly less than a half), when F-sentences have diminished by a little less than a third. All in all, it appears

the distribution of interrogative variants has shifted from the *ex situ* kinds to a more dominant *in situ*, over the course of half a century.
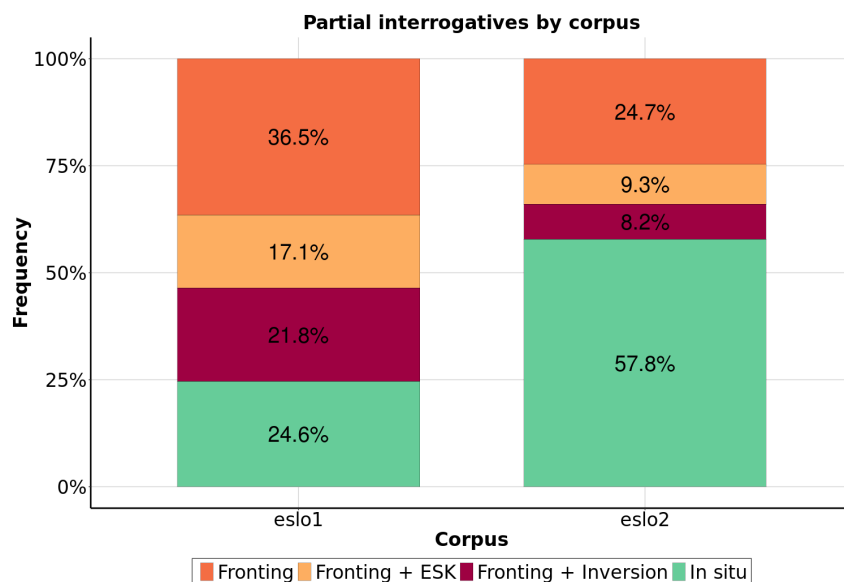


*figure 2*:  FPI type by corpus (%)

Bayesian modelling of this data refines the analysis of these broad evolutions[13]. A model with *FPI type* as the dependent variable was run, and with *corpus* as the main fixed predictor. A random factor was added: the *annotated file* where the FPI had been extracted from, to try and account for data variability. Convergence was reached and checked by making sure Rhat was equal to 1 for each parameter. The reference level of the dependent variable was *in situ*, IS, and the probability of each other levels being realized, given a possible influence of *corpus* and the random variables, was calculated against it. The model was run with 4 chains with 3000 iterations by chain (half of them for warming up the model, the other half for useful sampling).

With those parameters, an effect of *corpus* was found for all sentence type: F sentences ($\beta$= -1.36, 95%CrI=[-2.00,-0.76], P($\beta$)<0=1), Fesk sentences ($\beta$= -1.43, 95%CrI=[-2.10,-0.78], P($\beta$)<0=1), and FINV sentences ($\beta$= -2.09, 95%CrI=[-3.00,-1.23], P($\beta$)<0=1). P($\beta$) is here the

---

[13]    Here and thereafter, the general output of the models will be given in plain text, with the full results and parametrizations available on the OSF repository https://osf.io/ug8bt/ (see model 1 for this one).

probability that less F-/Fesk-/FINV- sentences were used in ESLO2 than in ESLO1, relatively to the evolution of IS-sentence uses from the first corpus to the second one. Figure 3 gives an illustration for the strength of the effects. 0 on the X-axis is the reference level for the comparison of distributions (evolution of IS sentences use between ESLO1 and ESLO2). Posterior distributions to the left of this threshold indicate how much the use of F/FESK/FINV sentences has decreased when compared to IS. A posterior distribution to the right of this threshold would have indicated a bigger increase in use than that of IS sentences.
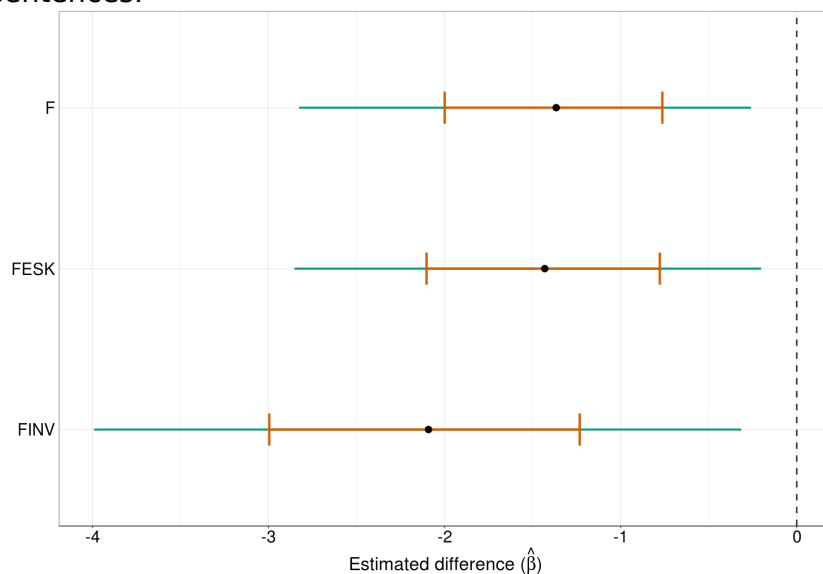


*figure 3*: Posterior distributions for FPI types evolution, relative to the evolution of IS sentences, with corpus as a main predictor (95% CrI)

This diachronic change is in line with other works studying the evolution of FPI use across time periods, with Larrivée (2016, 2019) also finding an evolution with less fronting and more *in situ* sentences in the modern era than in more ancient data, going back to the 15th century. In Larrivée (2019), this evolution of linguistic preferences is linked, based on corpus data also from the ESLO project, to the idea that *in situ* structures needed to be explicitly activated in the discourse in earlier times (by means of a declarative sentence resembling the interrogative one, see example (3) below), which is not the case anymore (4).

(3)  OW26 : Dans les jeux antiques euh ils se dopaient quand même avec des méthodes un peu bizarres mais

ch_PP6 : ils se dopaient comment ?
OW26 : ils prenaient euh des plantes (ESLO2_ENT_1026)

(4)    finalement tu trouves comment la vie à Orléans ?
(ESLO2_ENT_1021)

## 2.4.    The age factor: speakers from different age groups behave differently

Apart from the overall diachronic pattern, the frequencies of use are also age-related.

Experimental data (Thiberge, 2018) show a main effect of age in the appreciation of three variants of FPIs. In an acceptability judgment task, French adult native speakers (N=57, age 18-72, mean 28, median 33), read and rated interrogative sentences, with participants older than 30 years rating FINV consistently better than F or IS variants, as compared to participants less than 30 y.o.

Here, speakers were divided into two groups based on the INSEE age metadata available, 15-25 vs. 35-55 y.o. These age groups are the ones available in the INSEE data, and correspond to age groups below and above the 30-year threshold which appeared to be statistically meaningful in Thiberge (2018). The difference in frequency of use for each variant across the age groups doesn't necessarily seem very large at first (Fig. 4). However, there seems to be a slight decrease of IS uses for 35-55 speakers as compared to 15-25 speakers and a slight increase in F and FINV uses by 35-55 speakers.
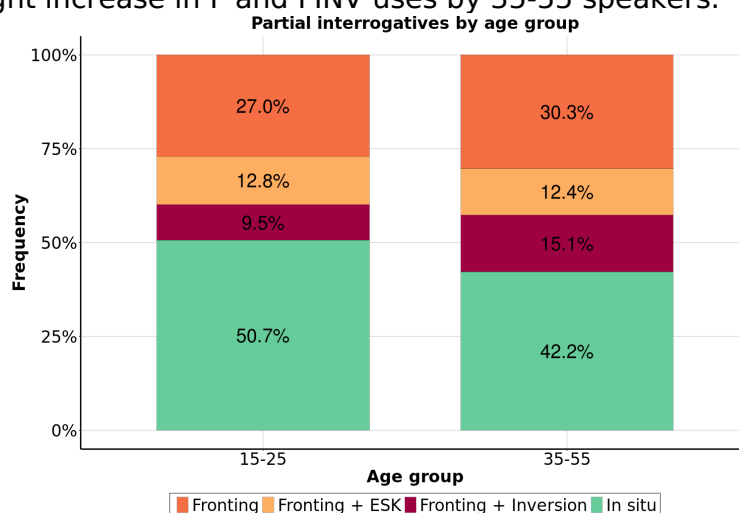


*figure 4*:  FPI type by age group (all tokens, %)

With a more in-depth take on the frequencies, things are different when both corpora are considered separately. Figure 5 illustrates how the variation of question type usage is both time-period-dependent and age-group-dependent. In the ESLO1 corpus, both the 15-25 and the 35-55 y.o. groups use a fronted structure for roughly 75% of the tokens. Note here the different weight of FINV uses between age groups, with a quarter of FPIs in FINV form for the 35-55 group, as opposed to 15-25 speakers who use this variant less than 15% of the time. In the ESLO2 corpus, the discrepancy between age groups is of a different magnitude. In the 2010's time period, both age groups use *in situ* (IS) structures for more than half of the tokens, but the 15-25 y.o. group makes even more use of those (~70%) than their elders (~55%), with a simultaneous big drop in all types of fronting uses. Comparatively to the decrease observed for younger speakers in F- (reduction by half) and Fesk- (reduction by two thirds) sentences, the phenomenon seem less important in the 35-55 year old group (reduction by a third).
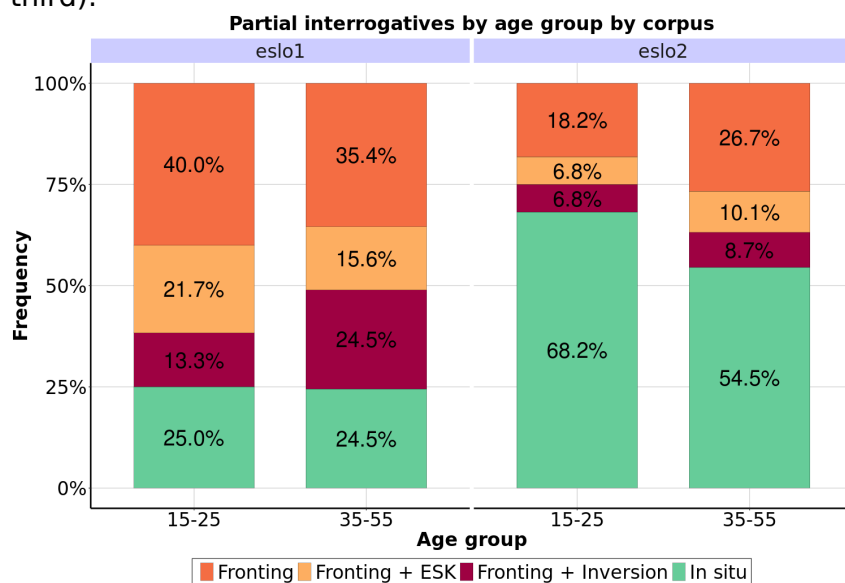


*figure 5*:  FPI type by age group by corpus (%)

Bayesian modeling of the data once again refines these observations.

A model for the surface influence of age was run, with *FPI type* as the dependent variable, and with *age group* as

the main fixed predictor1[14]. With these parameters, no main effect of *age* was found, for either F-/Fesk-/FINV-sentences[15]. As for a more in-depth look, the previous model was modified to add a main fixed predictor, *corpus*, in interaction with the *age group* variable. All other parameters were the same (random effect, chains, dependent variable and reference level); convergence was reached and controlled[16].

The main effect of *corpus* was confirmed for all three variants as was the overall absence of effect for *age* as a single predictive factor. Interactions between the two predictors led to somewhat nuanced results (Fig. 6). 0 on the X-axis is the reference level for the comparison of distributions (still, IS sentences).

---

[14]    A random factor with simple intercept was once again the *annotated file*. Convergence was reached and checked (Rhats = 1 for all parameters). The reference level of the variable of interest was still *in situ*, IS, and the probability of each other level being realized was calculated against it. The model was run with 4 chains with 3000 iterations by chain.

[15]    See model 2 in the OSF repository. With $P(\beta)$ being here the probability that more or less of a sentence type were used by the 35-55 y.o. group than by the 15-25 y.o. group, relatively to the difference of IS- uses between the two groups, for F: $\hat{\beta}$= -0.18, 95%CrI=[-0.97,0.57], $P(\beta)$<0=0.66; for Fesk: $\hat{\beta}$= -0.06, 95%CrI=[-0.94,0.81], $P(\beta)$<0=0.55); for FINV: $\hat{\beta}$= 0.41, 95%CrI=[-0.77,1.59], $P(\beta)$>0=0.75.

[16]    See model 3 in the OSF repository for a detailed presentation of the model. For clarity purposes, only the output for meaningful interactions is reported in the main text, but not the confirmatory output for both main effects.
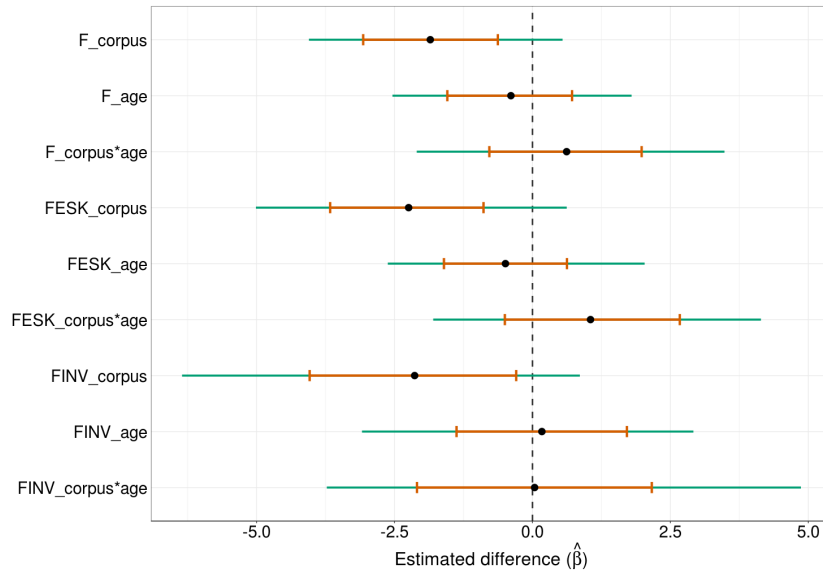
*figure 6*: Posterior distributions for FPI types evolution, relative to the evolution of IS sentences, with corpus and age group as main predictors (95% CrI)

While there seems to be no interaction of *corpus* and *age group* ("corpus*age"), for all types of fronting, the F- and Fesk- interactions are in fact meaningful, admittedly to a lesser degree (respectively, $\hat{\beta}= 0.61$, 95%CrI=[-0.78,1.98], P($\beta$)>0=0.81, and $\hat{\beta}= 1.05$, 95%CrI=[-0.50,2.67], P($\beta$)>0=0.91). P($\beta$) is here the probability that more F- / Fesk- sentences were used in ESLO 2 by the 35-55 group than by the 15-25 group, as compared to the difference of use between both age groups in ESLO 1, relatively to the between-groups evolution of IS-sentence uses from the first corpus to the second one. Put more simply, there's still an 81%/91% -chance that the 35-55 group are using more F-/FESK in ESLO2 than the 15-25 group because of a combination of the *corpus* and *age group* factors, when compared to their respective changes in IS uses.

In summary, from a diachronic perspective, corpus data provide us with two illustrations of how language preferences vary with time. First, in no more than fifty years – which is maybe *two* generations –, a social group (here, the people living in Orléans) can change their linguistic habits overall (shift from *ex situ* to *in situ* FPIs). Second, even when taking a more restricted time window for reference, it seems there is something like a "generation gap" with young speakers exhibiting different

preferences than their elders. This may provide further useful insight on how a given language change phenomenon can propagate within a linguistic community. The younger speakers in the 1960s used comparatively less frontings with inversion, in all their frontings, than the 35-55 age group of the same time window. Fifty years later, in the 2010s, the rarefaction of fronting with inversion is confirmed across both age groups, and now the main differences between them lie in their propensity to use an *in situ* construction or to "avoid" interrogative sentences with simple fronting.

This analysis needs to be nuanced, since the ESLO project is limited by design to a specific population (the inhabitants of the region surrounding the city of Orléans, who are probably not a perfect mirror for all francophone populations across the second half of the 20th century). The results are however consistent with other recent corpus studies (Adli, 2015; Hamlaoui, 2009) Another fact to take into account is the diversity of the linguistic interactions recorded, compiled and transcribed in the ESLO corpora.

## 3. A context-based corpus exploration of the use of FPIs variants

### 3.1. Overall relationship between context and interrogative structures

Sociolinguistic variation is a two-level phenomenon. It can be the reflection of a sociolect, i.e. *some groups of speakers make a higher use of one of the available variants (because of habits or because it signals their belonging to one or several social groups)*. This is the leading perspective under which sociolinguistic data has been analysed on the issue of French partial interrogatives (Pohl, 1965; Terry, 1970; Behnstedt, 1973; Ashby, 1977; Söll, 1982; Coveney, 1996). But variation can also be a tool, used by speakers to convey some social cues about themselves in a particular interaction, i.e. *the same speaker may use one or different linguistic variants in different social contexts, depending on what social persona they want to draw around themselves in their interlocutor's perception*. This dichotomy has fuelled an evolution within sociolinguistics throughout the second part of the 20th century, maybe mostly in the English literature; Eckert (2012) gives an account of these different approaches.

In such a refined perspective, corpus data must not only be analysed under a diachronic prism or with lenses that

divide speakers in multiple social groups (be it by age, socioeconomic origin or status, linguistic background, etc.). The context of interaction in which linguistic data was retrieved must also be taken into account, and in this respect the ESLO project is of big value. A great dimension of this corpus is the diverse methodology for data collection, ranging from supervised and hard-scripted interviews to recordings of families in fully spontaneous and unguided interactions at meal-time, with efforts made to use the same variety for both ESLO 1 and ESLO 2 and with the same protocol used for data collection in each context.

For the 617 annotated tokens that fit the elected criteria for statistical analysis, the list of the different interactive contexts where informants were recorded is as follows, with Table 3 summarizing the ESLO FPIs distribution over all contexts selected. It should be noted that all contexts retained for analysis are contexts where recorded speakers were speaking spontaneously, without any preparation or script (contrary to other kind of contexts available in the ESLO project, such as conference recordings for instance).

– INTERVIEW ("entretien"): interviews between a researcher and a linguistic informant, with discussions following a prepared and standardized outline (e.g. each informant was asked during the interview "how do you prepare an omelette?")
– SCHOOL ("école"): classes given by a teacher to pupils
– MOVIE ("cinéma"): short interviews by researchers outside of movie theaters, where they ask random persons a few questions, following a script
– ITINERARY ("itinéraire"): short interviews led by researchers in the streets of Orléans, where the researchers ask how to go to some places within the city boundaries, following a script
– MEAL-TIME ("repas"): home interactions at meal-time
– MEETING ("réunion"): business or public meetings
– PHONE ("téléphone"): phone calls to linguistic informants, administrations or shops
– 24H ("24h"): 24-hour-long recordings of someone's everyday interactions
– OTHER ("divers"): "opportunistic recordings" of various situations

*table 3*:    Distribution of FPIs by context (ESLO 1&2) and position of the Wh- element

|          | Question type | | | | TOTAL |
|----------|-----|------|------|-----|-------|
|          | f   | fesk | finv | is  |       |
| 24h      | -   | -    | 1    | 2   | 3     |
| movie    | 1   | -    | 1    | 14  | 16    |
| other    | 3   | -    | -    | 2   | 5     |
| school   | 39  | 7    | 19   | 65  | 130   |
| interview| 113 | 59   | 60   | 104 | 336   |
| meal-time| 1   | 1    | -    | 10  | 12    |
| meeting  | 17  | 6    | 3    | 68  | 94    |
| phone    | 8   | 4    | 1    | 7   | 20    |
| phone    | -   | -    | -    | 1   | 1     |
| Total    | 182 | 77   | 85   | 273 | 617   |

It appears here the data is quite unbalanced between contexts. Clearly, the distribution of FPIs over all different contexts is statistically biased (e.g. the PHONE context, where only one interrogative was kept for analysis), but even for the three main types of context (i.e. SCHOOL, INTERVIEW and MEAL-TIME), important phenomena appear when going into details.

Namely, the three contexts under consideration here do not elicit the same main question type (Fig. 7). While meal-time interactions elicit nearly 75% of IS partial interrogatives, interviews only elicit them in about a third of the utterances. By contrast, F sentences represent a third of all FPIs in interviews and at school, but less than 20% at meal-time. The importance of FINV forms in both school recordings and interviews (+/- 15%) and of Fesk sentences in interviews (more than 15%) is also worth noting.
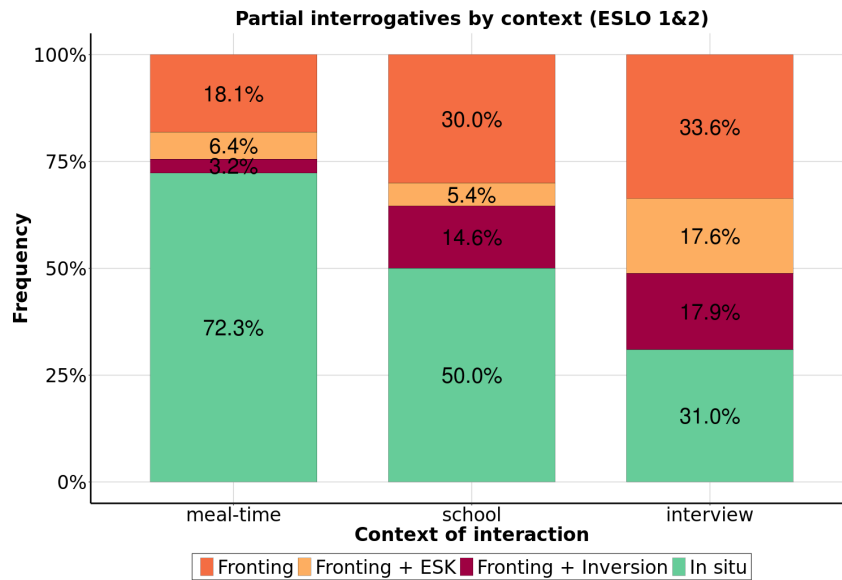
**Partial interrogatives by context (ESLO 1&2)**



*figure 7*:  FPI type by context of interaction in ESLO 1 & 2 (%)

A Bayesian model[17], run on the subset of data corresponding to these three specific contexts (N=560 out of 617 FPIs), refines this pattern. *FPI* was the dependent variable, and *context* was the main fixed predictor with three levels (interview/meal-time/school). A random factor with simple intercept was once again the *annotated file*. Convergence was reached and checked (Rhats = 1 for all parameters). The reference level of the variable of interest was still *in situ*, IS, and the probability of each other levels being realized was calculated against it. The reference level for the main predictor *context* was the SCHOOL context, because the proportions of IS used there seem to be some kind of middle ground (50%) between meal-time interactions (~75%) and interviews (~30%). The model was run with 4 chains with 3000 iterations by chain.

As compared to school interactions, interactions at meal times elicit less F- and FINV- sentences, relatively to the variability of IS uses between the two contexts (respectively, $\beta$=-1.29, 95%CrI=[-2.65,-0.07], P($\beta$)<0=0.98, and $\beta$= -2.40, 95%CrI=[-4.60,-0.39], P($\beta$)<0=0.99). P($\beta$) is here the probability that less F- / FINV- sentences were used during meals than at school. By contrast, as compared to school interactions, interviews elicit more Fesk sentences, relatively to the variability of IS uses between the two contexts ($\hat{\beta}$= 1.79,

---

17    See model 4 in the OSF repository.

95%CrI=[0.82,2.88],        P(β)>0=1).  P(β)  is  here  the
probability  that  more  Fesk-  sentences  were  used  during
interviews  than  at  school.  This    probability  is  also  quite
high  for  F-  sentences,  but  to  lesser  degree  (β=  0.57,
95%CrI=[-0.29,1.41],        P(β)>0=0.91).  Figure  8  is  a
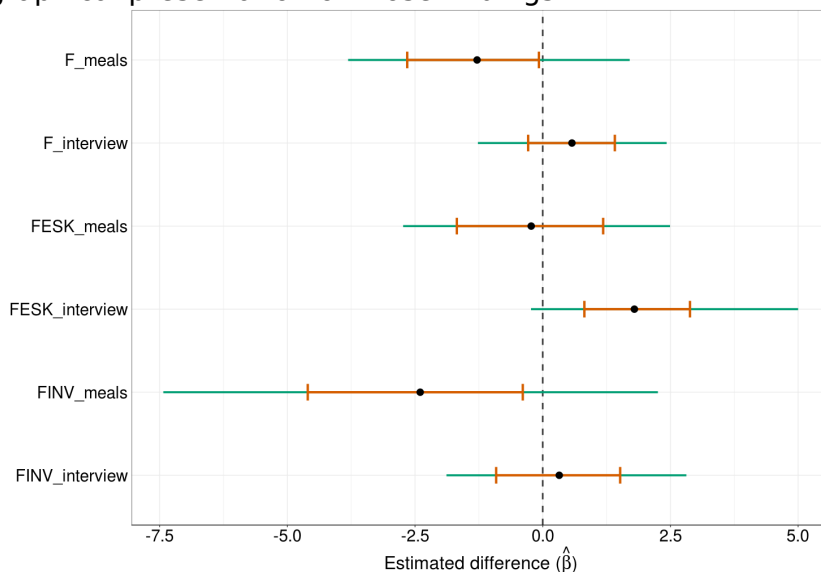graphical presentation of those findings.



*figure 8*: Posterior  distributions  for  FPI  types  variability  between
contexts, relative to the variability of IS uses (95% CrI)

As  mentioned  in  the  introduction,  even  though
pragmatic  factors  might  be  at  play  in  the  choice  of
interrogative  structures  (Hamlaoui  2009,  2010,  2011 ;
Déprez et al., 2013, a.o.), here it seems corpus data really
gives  ground  for  a  refined  socio-pragmatic  analysis  of
French  variation  regarding  partial  interrogatives.  On  the
one  hand,  yes,  French  native  speakers'  choices  in  the
matter could be explained by the social group they belong
to,  which  could  underscore  a  phenomenon  of  so-called
"pure optionality": the choice would be self-conscious, and
just  a  means  to  express  the  speakers'  identity  as
members of a social group (defined by age, for instance).
On  the  other  hand,  it  could  be  much  more  nuanced,  and
the  qualities,  or  perceived  qualities,  of  a  context  of
interaction  could  influence  the  way  speakers  behave  in
this  interaction.  Here,  it  appears  meal-time  interactions
give  rise  to  a  prevalent  use  of  IS  structures  and  a  much
reduced use of FINV structures, as opposed to interviews
between  an  informant  and  a  researcher,  which  give  rise  to
more fronting, with or without subject-verb inversion, with

or without "est-ce que" insertion. School recordings, where the interrogatives are spontaneously produced mainly by the teacher talking to children a classroom, would fall somewhat in the middle.

Those three contexts differ in many aspects: the interlocutor is not the same (children vs. adults, external researcher vs. familiar figures) but also the setting is very different (public space vs. home). The combination of those differences could render a somewhat hierarchical pattern, with interviews being perceived as the most formal – or "not-so-colloquial" – of the three contexts, meal-time interactions being of a much more colloquial nature, and school interactions being some kind of formality-wise middle-ground. The precise origin of this difference between contexts would need further exploration, and some authors associate this (in)formality with an opposition between "controlled" vs. "non-controlled speech" (Zribi-Hetz, 2011) in a diglossic approach to French.

Those results are coherent with Thiberge (2018), where it was found experimentally that people using FINV interrogatives were perceived as behaving less "relaxed", and thus acting more "formally", than the users of IS and F sentences, but also that the use of FINV forms by someone entailed a higher probability of him/her being rich, a frequent reader, and educated. But there is no need to try and appear rich, a frequent reader or educated, when interacting with family members at meal time. The balance of powers between participants is however different when they are not from the same world, as in interviews where a researcher comes into a home to ask someone about their linguistic habits, or when there is a hierarchical relationship of sorts, as in school between a teacher and their pupils. In those two particular cases, the concept of social – and public – persona (Ochs, 1992) comes into play. This is in line with recent sociolinguistic works where people adopt different linguistic behaviours depending on the context of interaction, the audience, and their communicative strategies (see Labov (2012) for an analysis of President Obama's production of the -ing/in' verbal variants during informal barbecues, political interviews or formal speeches).

## 3.2.  Combining the diachronic and social group approaches: data from interviews

This context-based analysis of linguistic interactions and of sociolinguistic variation can be combined with the micro-diachronic perspective adopted in the previous section. It is quite possible to look at whether – and how –

social groups (15-25 y.o. vs. 35-55 y.o.) and time periods (1960s vs. 2010s) interact with the use of all three variants within a specific context of interaction. Only interview transcriptions can be compared across time periods, since the partial interrogatives from ESLO1 meal-time interactions were only produced by speakers from the 35-55 y.o. group.

Interview data show a superficial difference across age groups in the FPI construction: when analysing aggregated data from both the ESLO1 and ESLO2 corpora, the 15-25 y.o. group use more *in situ* constructions in this context than the 35-55 y.o. group overall, with the latter exhibiting the same slightly broader use of FINV structures as seen before (subsection 2.4). Figure 9 gives an illustration of this 'generation gap', while Bayesian modelling of the data[18] is only tentatively conclusive, with a mere 74% chance that more FINV sentences are used by 35-55 year-old speakers than by 15-25 years old speakers ($\hat{\beta}$= 0.41, 95%CrI=[-0.91,1.71], P($\beta$)>0=0.74).

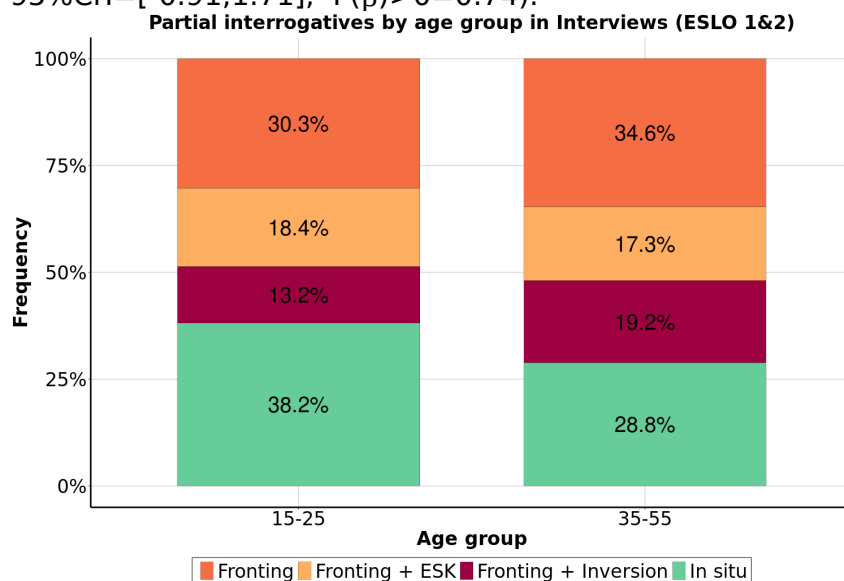**Partial interrogatives by age group in Interviews (ESLO 1&2)**



*figure 9*:  FPI type by age group in both ESLO 1 & 2 (interviews, %)

Second, this generation gap is strikingly not the same across time periods. In the same pattern observed previously, there is less of a difference between age

---

[18]     See model 5 in the OSF repository, which was run with the exact same parameters as model 2, but on a smaller set of the data (N= 336 out of 617 FPIs). Try-models with more iterations by chains lead to no improvement whatsoever.

groups in the ESLO1 data than in the tokens extracted from ESLO2 (Fig. 10). All speakers from the 1960s interviews were using far less *in situ* constructions than frontings. There is a big difference of proportions between all three types of fronting (F vs. Fesk vs FINV), but both groups produce, in this specific context of interaction, only around 20% of *in situ* interrogatives. On the other hand, in the data from the 2010s, the gap seems to materialize differently, in the same specific context of interaction. While there is now a generally reduced use of FINV structures, 35-55 y.o. speakers use *in situ* constructions less than 50% of the time, while 15-25 y.o. speakers use them nearly in 60% of their productions. Moreover, in a pattern quite similar to what was observed before in the un-contexualized data (Fig. 5), F sentences but also Fesk sentences seem to be stronger alternatives to *in situ* FPIs for participants in the 35-55 group than for the 15-25 group.
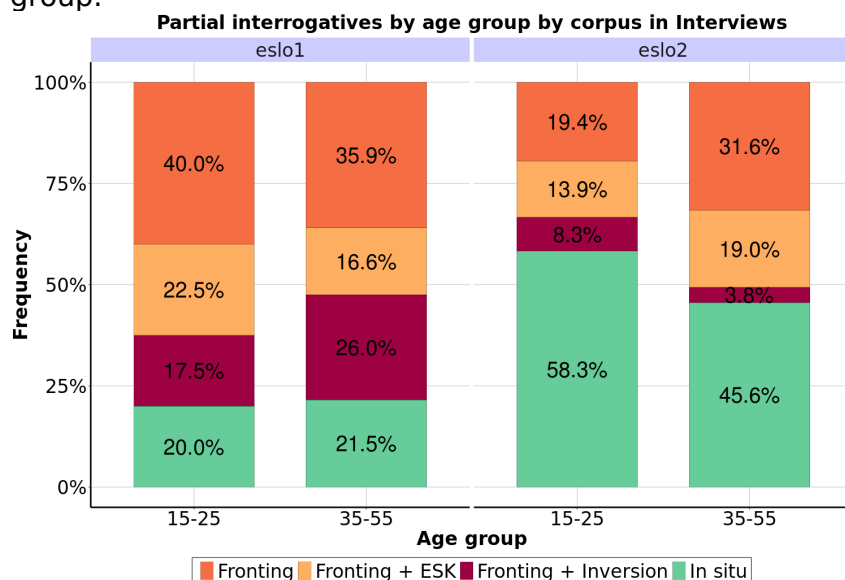


*figure 10*: FPI type by age group by corpus (interviews, %)

A Bayesian take on this subset of extracted FPIs1[19] return the same pattern of results as on the whole corpus. The main effect of *corpus* is confirmed, and no overall

---

[19] See model 6 in the OSF repository for details. Again, the model was run with the exact same parameters as model 3, but on a smaller set of the data (N= 336 out of 617 FPIs). For clarity purposes, and as the results are absolutely comparable to those of model 3, the numeric values were not put in plain text.

effect of *age group* can be confirmed statistically. However, interactions between *corpus* and *age group* are meaningful to some extent, with a reasonable probability of this combination of predictors leading to a more frequent use of F- and Fesk- by 35-55 year old speakers than by 15-25 year old speakers, in ESLO 2, relatively to the difference of their IS uses from ESLO1 to ESLO2 (82% for F- sentences 88% for Fesk- sentences).

### 3.3.    Combining social group and contextual analyses: contrasts between interview and meal-time data (ESLO2)

Data from the ESLO2 interview subset can finally be compared to the data from meal-time interactions (only available for both age groups in the ESLO2 corpus). Figure 11 show how the two age groups compare on their respective uses of each FPI variant, in both interviews and meal-time interactions. During meals, *in situ* interrogatives are produced in an equally high proportion of ~70% or more of the time by both age groups. Numerically, the 35-55 y.o. group actually uses more *in situ* FPIs when compared to the younger group. This strikingly contrasts with data seen in the previous subsection for interviews, where *in situ* structures were not totalling half of the tokens for the 35-55 group (for less than 60% of uses for the 15-25 group).
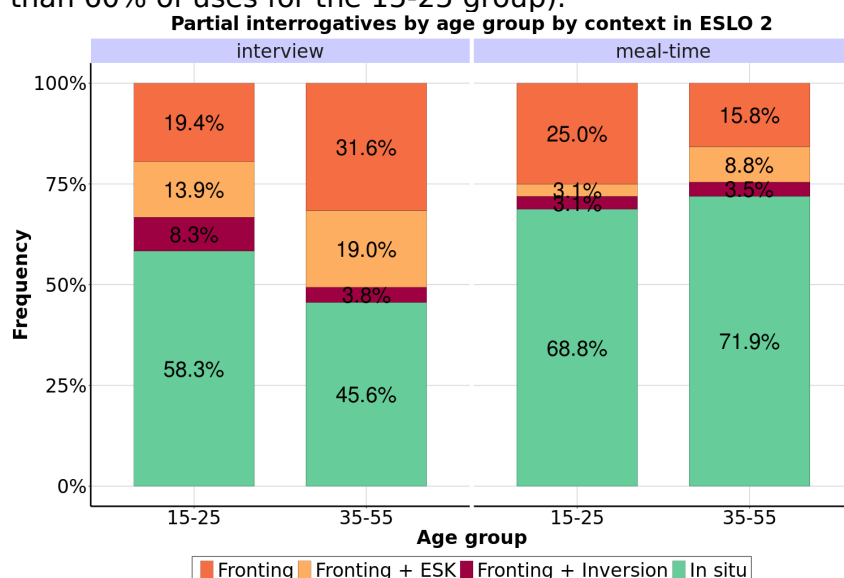
**Partial interrogatives by age group by context in ESLO 2**

| | interview | | meal-time | |
|---|---|---|---|---|
| | 15-25 | 35-55 | 15-25 | 35-55 |
| Fronting | 19.4% | 31.6% | 25.0% | 15.8% |
| Fronting + ESK | 13.9% | 19.0% | | 8.8% |
| Fronting + Inversion | 8.3% | 3.8% | 3.1% 3.1% | 3.5% |
| In situ | 58.3% | 45.6% | 68.8% | 71.9% |

*figure 11*: FPI type by age group by context, in ESLO2 (%)

Bayesian modeling of the data follows the trends visible in raw frequencies[20]. *FPI type* was, as always, the dependent variable, and both *context* and *age group* were used as the main fixed predictors. The *annotated file* was still added as a random factor. Convergence was reached and checked (Rhats =1 for all parameters). The reference level of the dependent variable was IS, and all computations were ran by using 'interview' as the reference level of the *context* parameter. The model was again run with 4 chains with 3000 iterations by chain.

Though nothing jumps out when considering 95% credibility intervals, an isolated main effect of context was close to being salient for Fesk- sentences ($\beta$= -2.25, 95%CrI=[-5.97,0.63], P($\beta$)<0=0.94), with P($\beta$) being the probability that less Fesk- sentences were used globally in ESLO 2 for both age groups during meals, as compared to interviews, relatively to the global evolution of IS uses. An interaction between the two main predictors, *context\*age*, was found to be quite probable for F- sentences ($\beta$=-1.49, 95%CrI=[-3.70,0.70], P($\beta$)<0=0.91). There, P($\beta$) would be the probability that less F- sentences were used by the 35-55 year old speakers during meal-times than by speakers aged 15-25, relatively to their change of uses in IS-sentences between the two contexts. The change of linguistic behaviour in the 35-55 year old population (reduction of F- and augmentation of IS- production) from interviews to meals is then of a much more important magnitude than the change (also visible) for the younger group.

Here again, the concept of social or public persona (Ochs, 1992) can be useful. The difference in perceived level of formality between interview contexts and meal-time interactions can explain this speakers' behaviour. In a context where formality is irrelevant and the expectation of speaking "good French" is not applicable, it doesn't seem surprising that both age groups would mainly use the IS structure, less associated with marked social cues such as richness, education or frequent reading (Thiberge 2018). On the other hand, in an INTERVIEW context, where a researcher from the outside is asking questions or asking for someone to discuss a variety of topics for research purposes, even in a voluntarily familiar setting (the interviews were conducted at home), it may be inevitable that informants would be careful of how "well" they speak. If what someone says and the way they say it affects how they are perceived, then it is only logical that someone would adapt their linguistic productions depending on the

---

[20] See model 7 in the OSF repository.

context, the other persons they are speaking to, and what they want to convey about themself.

But this only explains the overall difference of proportions between *in situ* structures and the variety of frontings available to French speakers. To explain the discrepancy between age groups in their use of frontings vs. *in situ* constructions in ESLO2 interviews, one could postulate that speakers aged 35-55 behave differently with regards to formality than speakers aged 15-25[21]. More precisely, either speakers from the 35-55 group feel the need to speak more formally than speakers aged 15-25 or it is the other way around and speakers aged 15-25 are less attentive to their level of formality than speakers aged 35-55. Given the fact that neither group makes a substantially broader use of the FINV structure (mostly associated with marked social cues such as richness, education and frequent reading, as per Thiberge(2018); and see Adli (2015) or Hamlaoui (2009) for other recent corpus studies where FINV uses are very limited if not absent), one could think the second hypothesis is more plausible and the younger speakers feel less compelled to obey sociolinguistic norms, which would need to be tested in a more systematic way through linguistic experiments.

It could also be the case that the sociolinguistic norms weighing on the younger speakers have come to be different than those of their elders through language

---

[21]    Experimental data yet to be published actually give an illustration of this. When native speakers are asked whether they think an interrogative variant is "acceptable French", they tend to answer differently when the variant is preceded with a formal context or with an informal context. Participants aged more than 30 y.o. – who in Thiberge (2018) were the only ones significantly differing in their acceptability judgments of FINV and IS/F structures – exhibit yet a higher sensitivity than  participants aged less than 30 y.o.. See Thiberge & Hemforth (2019), poster at the 2018 AMLaP conference available online: http://www.llf.cnrs.fr/sites/llf.cnrs.fr/files/biblio//180907%20AMLaP%20poster.pdf Thiberge (2020) presents further experimental work along this line, with acceptability judgment tasks where speakers aged between 40-60 y.o. give worse ratings to IS sentences in formal contexts than other participants, which. This effect resembles the notion of "age-grading" (Wagner, 2012), according to which people of this age group are the ones most confronted to linguistic norms, for example through profesional interactions, and are most likely to have strong judgments on which linguistic variants are acceptable in formal contexts as opposed to other variants.

change processes. Either way, this could again be a further illustration of how linguistic change is a dynamic phenomenon, rooted in the change of linguistic habits from one social group (here, again, defined by age).

It could also be the case that the sociolinguistic norms weighing on the younger speakers have come to be different than those of their elders through language change processes. Either way, this could again be a further illustration of how linguistic change is a dynamic phenomenon, rooted in the change of linguistic habits from one social group (here, again, defined by age).

## Summary, conclusion and further work

Spontaneous French data from the ESLO corpora, taken in its diachronic dimension and the variety of interactive contexts it provides, allowed for a fine-grained analysis of French partial interrogatives. Data showed an evolution in the use of syntactic structure of those sentences, most visibly from a general pattern of *ex situ* constructions to a more generalized *in situ* realization of the Wh- element. Different patterns were also seen between contexts, in the sense that less informal contexts (interviews with a researcher) or more public interactions (school) elicited comparatively less *in situ* productions, as opposed to interactions in a more familiar environment like meal-time interactions with family members. A third facet of interest of the ESLO project was exploited, namely the sociolinguistic metadata of its speakers, which led to a multifactorial approach to the alternation observed in modern French partial interrogatives.

Not only is the FPI type dependent on the context of interaction and in the time window of analysis (mid-20[th] century of dawn of the 21[st] century), but the age-group to which a speaker belongs will also weigh on his or her choices with regards to the production of a partial interrogative. A similar linguistic behaviour in informal contexts with a small familiar audience (family) for both age groups considered (15-25 y.o. vs. 35-55 y.o.) but more contrasted behaviours in more open settings (interviews) is taken to be an illustration of how group identity shape, in proportion, individual communicative strategies. In a broader perspective, these findings argue for a more nuanced analyses even of conversational data: dialogue is a social activity and social strategies must be taken into account for a full analysis of syntactic phenomena. This general finding is in line with recent experimental 3rd-wave sociolinguistic works, in the experimental field for example.

Several open questions remain and should be the object of further investigations:

1°/ Four kinds of FPIs where contrasted here, F, Fesk, FINV and IS. They have all been treated as independent from one another, but the question of the mutual relationships between all fours variants needs be more thoroughly examined, beyond the simple opposition between *ex situ* and *in situ* constructions.

2°/ Some idiomatic constructions (*comment dirais-je?* – "how could I say this") were not specifically annotated as such and might be of interest.

3°/ The point of mirror-questions, that is questions that are answering a question by asking a new one with the same structure as the first one, were not analysed because of the automatic extraction of the data. Further analysis on this aspect might reveal new facets to the specific phenomenon of *variation in FPIs*, as well as bring forth new insight on the spontaneous co-building of (socio-)linguistic interactions.

4°/ Because the protocols for data collection were the same for the ESLO1 and the ESLO2 sub-corpora in similar contexts and across age groups, the quantitative approach and the statistical inferences from Bayesian modelling were taken as sufficient evidence for diachronic, diastratic and diaphasic parameters influences on the uses of FPIs. The balance of influence between these social parameters and more linguistic factors such as pragmatic considerations needs further investigation. This could be done by a reannotation of the data, oriented towards a discourse/situational analysis of FPIs production, and with more experimental approaches to the variation.

## Access to data

The materials (all extracted interrogatives used for analysis, and the Rstudio script for the bayesian analysis conducted) are or will soon be made available online at https://osf.io/ug8bt/ .

## Acknowledgments

## Bibliography

Analyse et traitement informatique de la langue française - UMR 7118 (ATILF), Institut de l'information scientifique et

technique - CNRS UPS76 (INIST), Laboratoire d'informatique de Paris Nord - UMR 7030 (LIPN), 2012, *PERCEO : un Projet d'Etiqueteur Robuste pour l'Ecrit et pour l'Oral* [Corpus]. ORTOLANG (Open Resources and TOols for LANGuage) - www.ortolang.fr, https://hdl.handle.net/11403/perceo/v1.

ASHBY W., 1977, Interrogative forms in Parisian French, *Semania* 4, p. 35–52.

ADLI A., 2015, What you like is not what you do : Acceptability and frequency in syntactic variation, in A. Adli, M. García García & G. Kaufmann (eds), *Variation in Language: System- and Usage-based Approaches*, Berlin, München, Boston, De Gruyter, p. 173–200. https://doi.org/10.1515/9783110346855

BAIÃO, V. L. & Lobo M., 2014, Aquisição de interrogativas preposicionadas no português europeu, in *Textos Selecionados, XXIX Encontro Nacional da Associação Portuguesa de Linguistica*, Porto, APL, p. 57–70.

BARRAS, C., Geoffrois, E., Wu, Z. & Liberman, M., 2001, Transcriber: Development and Use of a Tool for Assisting Speech Corpora Production, *Speech Communication* 33, p. 5–22.

BEHNSTEDT P., 1973, *Viens-tu ? Est-ce que tu viens ? Tu viens ? Formen und Strukturen des direkten Fragesatzes im Französischen*, Tübingen, Narr.

BENZITOUN, C., Fort, K. & Sagot, B., 2012, TCOF-POS : un corpus libre de français parlé annoté en morphosyntaxe, in *Actes de la conférence conjointe JEP-TALN-RECITAL 2012*, volume 2: TALN, Grenoble, p. 99–112.

BEYSSADE C., 2006, La structure de l'information dans les questions : quelques remarques sur la diversité des formes interrogatives en français, *LINX*, *Revue des linguistes de l'université Paris X Nanterre* 55, p. 173–193.

BOECKX C., 1999, Decomposing French questions, *University of Pennsylvania Working Papers in Linguistics* 6, p. 69–80.

BUERKNER P-C., 2017, brms: An R Package for Bayesian Multilevel Models Using Stan, *Journal of Statistical Software* 80(1), p. 1-28. https://doi.org/10.18637/jss.v080.i01

BUERKNER P-C., 2018, Advanced Bayesian Multilevel Modeling with the R Package brms, *The R Journal* 10(1), p. 395–411. https://doi.org/10.32614/RJ-2018-017

BURNETT H., 2017, Sociolinguistic Interaction and Identity Construction: The View from Game-Theoretic Pragmatics, *Journal of Sociolinguistics* 22(1), p. 238–271.

CARPENTER, B., Gelman, A., Hoffman, M., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M., Guo, J., Li, P. & Riddell, A., 2017, Stan: A Probabilistic Programming Language, *Journal of Statistical Software* 76(1), p. 1–32. https://doi.org/10.18637/jss.v076.i01

CLECH-DARBON, A., Rebuschi, G. & Rialland, A., 1999, Are there cleft-sentences in French?, in G. Rebuschi & L. Tuller (eds), *The grammar of focus,* Amsterdam, John Benjamins, p. 83–118.

COVENEY A., 1995, The use of the QU-final interrogative structure in spoken French, *Journal of French Language Studies* 5, p. 143–171.

COVENEY A., 1996, *Variability in spoken french: a sociolinguistic study of interrogation and negation*, Exeter, Elm Bank Publication.

COVENEY A., 1997, L'approche variationniste et la description de la grammaire du français: le cas des interrogatives, *Langue française* 115(1), p. 88–100.

COVENEY A., 2011, L'interrogation directe, *Travaux de linguistique* 63, p. 112–145.

DELAVEAU A., to appear, Les phrases interrogatives, in A. Abeillé & D. Godard (eds), *La grande grammaire du français*, Actes sud.

DÉPREZ, V., Syrett, K. & Shigeto, K., 2013, The interaction of syntax, prosody, and discourse in licensing French wh-in-situ questions, *Lingua* 124, p. 4–19. https://doi.org/10.1016/j.lingua.2012.03.002

ECKERT P., 2012, Three waves of variation study: The emergence of meaning in the study of sociolinguistic variation, *Annual review of Anthropology* 41, p. 87–100.

ENGDAHL E., 2006, Information packaging in questions, *Empirical issues in syntax and semantics* 6, p. 93–111.

ESKHOL-TARAVELLA, I., Baude, O., Maurel, D., Hriba, L., Dugua, C. & Tellier, I., 2012, Un grand corpus oral « disponible » : le corpus d'Orléans 1968-2012, in *Ressources linguistiques libres, TAL* 52, p. 17–46.

GOODMAN, N. D. & Lassiter, D., 2015, Probabilistic Semantics and Pragmatics: Uncertainty in Language and Thought, in C. Fox & S. Lappin (eds), *The handbook of contemporary semantic theory*, p. 655–686.

HAMANN C., 2006, Speculations about early syntax: the production of wh-questions by normally developing French children and French children with SLI, *Catalan Journal of Linguistics* 5, p. 143–189.

HAMLAOUI F., 2008, Focus, contrast, and the syntax-phonology interface: the case of French cleft-sentences, *Current Issues in Unity and Diversity of Languages, Collection of the papers selected from the 18th International Congress of Linguists (CIL18)*, Seoul, Linguistic Society of Korea.

HAMLAOUI F., 2009, *La focalisation à l'interface de la syntaxe et de la phonologie: le cas du français dans une perspective typologique*, Ph.D dissertation, Université Paris III Sorbonne Nouvelle.

HAMLAOUI F., 2010, A prosodic study of wh-questions in French natural discourse, in K. Clackson, Z. Absi, M. Ogawa, M. Ono, C. Patterson & V. Villafaña (eds), *Proceedings of the LangUE 2009*, Essex, University of Essex, p. 27–38.

HAMLAOUI F., 2011, On the role of phonology and discourse in Francilian French wh-questions, *Journal of Linguistics* 47, p. 129–162.

HEIDEN, S., Mague, J-P. & Pincemin, B., 2010, TXM : Une plateforme logicielle open-source pour la textométrie – conception et développement, in I. C. Sergio Bolasco (ed), *Proceedings of 10th International Conference on the Statistical Analysis of Textual Data - JADT 2010*, vol. 2, Roma, Edizioni Universitarie di Lettere Economia Diritto, p. 1021–1032.

HULK A., 1996, The Syntax of Wh-Questions in Child French, *Amsterdam Series in Child Language Development* 5, p. 129–172.

JAKUBOWICZ, C. & Strik, N., 2008, Scope-marking strategies in the acquisition of long distance wh-questions in French and Dutch, *Language and Speech* 51, p. 101–132.

JAKUBOWICZ C., 2011, Measuring derivational complexity: New evidence from typically developing and SLI learners of L1 French, *Lingua* 121, p. 339–351.

KRIFKA M., 2007, Basic notions of information structure, in C. Féry, G. Fanselow & M. Krifka (eds), *Interdisciplinary Studies of Information Structure* 6, Potsdam, Universitätsverlag Potsdam, p. 13–56.

LABOV W., 2012, *Dialect diversity in America : The politics of language change*, Charlottesville, University of Virginia Press.

LARRIVÉE P., 2016, Les interrogatives in−situ sont−elles pragmatiquement marquées en français vernaculaire? investigation synchronique et historique, in *4e colloque Repenser l'histoire du français − 7−8 avril 2016*, Münche, Ludwig Maximilians Universität München.

LARRIVÉE P., 2019, Historical pragmatics, explicit activation and whin situ in french, *Romance Languages and Linguistic Theory* 15, p. 114–132.

LEWIS D., 1969, *Convention. A Philosophical Study*, Harvard, Mass, Harvard University Press.

Laboratoire Ligérien de Linguistique - UMR 7270 (LLL), 2017, ESLO [Corpus], ORTOLANG (Open Resources and TOols for LANGuage) - www.ortolang.fr, https://hdl.handle.net/11403/eslo/v1.

LIÉGEOIS, L., Etienne, C., Parisse, C., Benzitoun, C. & Chanard, C., 2015, Using the TEI as a pivot format for oral and multimodal language corpora, *Text Encoding Initiative Conference and Member's meeting 2015*, Lyon, France.

Modèles, Dynamiques, Corpus - UMR 7114 (MoDyCo), 2016, Teicorpo [Outil]. ORTOLANG (Open Resources and TOols for LANGuage) - www.ortolang.fr, https://hdl.handle.net/11403/teicorpo/v1.

OCHS E., 1992, Indexing gender, in A. Duranti & C. Goodwin (eds), *Rethinking Context: Language as an Interactive Phemonenon*, Cambridge, Cambridge University Press, p. 335–358.

PARISSE, C. & LeNormand M-T., 2006, Une méthode pour évaluer la production du langage spontané chez l'enfant de 2 à 4 ans, *Glossa* 97, p. 20–41.

POHL J., 1965, Observations sur les formes d'interrogation dans la langue parlée et dans la langue écrite non littéraire, in *Actes du Xe Congrès International de Linguistique et de Philologie Romanes,* Tome 2, Paris, Klincksieck.

PRÉVOST P., 2009, *The acquisition of French: The development of inflectional morphology and syntax in L1 acquisition, bilingualism and L2 acquisition*, Amsterdam, John Benjamins.

QUILLARD V., 2001, La diversité des formes interrogatives : comment l'interpréter ?, *Langage et société* 95(1), p. 57–72.

R Development Core Team, 2009, *R: a language and environment for statistical computing*.

SCHMID H., 1994, Probabilistic Part-of-Speech Tagging Using Decision Trees, in *Proceedings of International Conference on New Methods in Language Processing*, Manchester.

SCHMID H., 1995, Improvements in Part-of-Speech Tagging with an Application to German, in *Proceedings of the ACL SIGDAT-Workshop*, Dublin, p. 47–50.

SOARES C., 2003, The C-domain and the acquisition of European Portuguese: the case of wh-questions, *Probus* 15, p. 147–176.

SOARES C., 2006, *La Syntaxe de la Périphérie Gauche en Portugais Européen et son Acquisition*, Ph.D dissertation, University of Paris 8.

SÖLL L., 1982, L'interrogation directe dans un corpus en langage enfantin, in F-J. Haussman (ed), *Études de grammaire française descriptive*, Heideberg, Groos.

SORENSEN, T., Hohenstein, S. & Vasishth, S., 2016, Bayesian linear mixed models using Stan: A tutorial for psychologists, linguists, and cognitive scientists, *The Quantitative Methods for Psychology* 12(3), p. 175–200. https://doi.org/10.20982/tqmp.12.3.p175

STRIK N., 2006, L'acquisition des phrases interrogatives chez les enfants francophones, *Psychologie Française* 52, p. 27–39.

STRIK N., 2008, *Syntaxe et acquisition des phrases interrogative en français et en néerlandais: une étude comparative*, Ph.D dissertation, University of Paris 8.

TERRY R. M., 1970, *Contemporary French interrogative structures*, Montréal, Sherbrooke, Éd. Cosmos.

THIBERGE G., 2018, Position du syntagme Wh- en français : réelle optionnalité ou biais sociolinguistique ?, *ELIS, Échanges de Linguistique en Sorbonne* 5, p. 64–91.

THIBERGE G., 2020, *Acquisition et maîtrise des interrogatives partielles en français : La variation comme outil interactionnel*, Ph.D dissertation, Université de Paris.

WAGNER S. E., 2012, Age grading in sociolinguistic theory, *Language and Linguistics Compass* 6, p. 371–382.

YEE, T. W., 2015, *Vector Generalized Linear and Additive Models: With an Implementation in R*, New York, Springer.

YEE, T. W. & WILD, C. J., 1996, Vector generalized additive models, *Journal of Royal Statistics Society* B58(3), p. 481–493.

ZRIBI-HERTZ A., 2011, Pour un modèle diglossique de description du français: quelques implications théoriques, didactiques et méthodologiques, *Journal of French Language Studies* 21.2, p. 231–256.

ZUCKERMAN, S. & Hulk, A., 2001, Acquiring optionality in French wh-questions: An experimental study, *Revue québécoise de linguistique* 30(2), p. 71–97.