

Improving the Science of Annotation for Natural Language Processing

KYLIE ANGLIN^{*1}, ARIELLE BOGUSLAV², AND TODD HALL²

¹*Department of Educational Psychology, University of Connecticut, Storrs, CT, USA*

²*Department of Education Leadership Foundations and Policy, University of Virginia,
Charlottesville, VA, 22902*

Abstract

Text classification has allowed researchers to analyze natural language data at a previously impossible scale. However, a text classifier is only as valid as the the annotations on which it was trained. Further, the cost of training a classifier depends on annotators' ability to quickly and accurately apply the coding scheme to each text. Thus, researchers need guidance on how to generate training data with optimal efficiency and accuracy. To this end, this study proposes the single-case study design as a feasible and causally-valid research design for empirical decision-making in annotation projects. The key strength of the design is its ability to generate causal evidence with as few as one annotator. In this paper, we demonstrate the application of the single-case study in an applied experiment and argue that future researchers should incorporate the design into the pilot stage of annotation projects so that, over time, a causally-valid body of knowledge regarding the best annotation techniques is built.

Keywords *annotation; coding; NLP; text classification; single-case study.*

1 Introduction

Text documents provide a rich source of data for linguists and social scientists alike. As these researchers bring their analyses to scale, text classification is playing an increasingly important role across many research domains. Without text classification, traditional

*Corresponding author. Email: kylie.anglin@uconn.edu

1 approaches to analyzing natural language data rely on highly trained personnel to read
2 and annotate each document of interest. First, researchers use theory to define a *coding*
3 *scheme*: a systematic framework for labeling each document that is tailored to the specific
4 research question (Shaffer and Ruis, 2021). Second, researchers either apply those labels
5 to the documents themselves or hire annotators to read and code each text. In addition
6 to participating in extensive training processes, hired annotators are often required to
7 have relevant professional and educational experiences that relate to the project’s specific
8 research area. The coding process is therefore both time and resource-intensive, limiting
9 the number of documents that can be included in any given study.

10 Text classification methods, on the other hand, only require hand-labeling text for
11 a subset of the documents comprising the training data. These data are then used to
12 train an algorithm to automatically apply the coding scheme to the remaining documents
13 in the corpus. Notably, once the text classification algorithm has been trained, it can
14 be applied again and again to additional documents at negligible cost. This feature
15 makes text classification a powerful and efficient analytic tool when study populations
16 are large and the amount of text data is voluminous. However, the validity and cost of
17 text classification depends on annotators’ ability to apply the coding scheme accurately
18 and efficiently to each text.

19 Unfortunately, researchers can currently find only limited guidance on how to pro-
20 duce valid and efficient hand-labelled training data. While qualitative social scientists
21 have devoted substantial attention to methods of producing valid human codes (Creswell
22 and Miller, 2000), this field has not traditionally needed to be concerned with the effi-
23 ciency of the coding procedure, nor whether the labelled text, and the complex social
24 constructs represented by those labels, are appropriate for automatic text classification.
25 For more specific advice on producing training data for machine learning, researchers
26 may instead turn to the field of computational linguistics where researchers have recently
27 begun to build a “science of annotation”, advising researchers on the best methods of
28 producing training data (Hovy and Lavid, 2010). However, this science of annotation
29 is still in a nascent stage. Despite growing calls for researchers to document the origins

1 and appropriate uses of training data in data statements or data sheets (Geburu et al.,
2 2021; Bender and Friedman, 2018), many papers today still fail to report key infor-
3 mation on how their training data were obtained (Geiger et al., 2020). Further, while
4 researchers can find insightful and practical recommendations (Hovy and Lavid, 2010; Ide
5 and Pustejovsky, 2017; Pustejovsky and Stubbs, 2012), there are still few studies that
6 have empirically tested the best methods of annotation. A few key exceptions include
7 research on the potential influence of annotator characteristics (Alyuz et al., 2021; Snow
8 et al., 2008), of iterative consensus building among annotators (D’Mello, 2016), and of
9 pre-annotation (Lingren et al., 2014), as well as several reviews of annotation software
10 (Dipper et al., 2004; Neves and Ševa, 2021; Neves and Leser, 2014).

11 Currently, empirical tests of annotation techniques commonly take one of two forms;
12 both of which create challenges for identifying the causal effect of particular annotation
13 methods. In one common approach, researchers may use a pre-post design where anno-
14 tators use one method of annotation followed by an alternative method. Performance
15 statistics are then compared across the time points. This design is straight-forward
16 but presents severe challenges for causal inference. Namely, it is impossible to decipher
17 whether changes in annotator performance are due to the new method of annotation,
18 or due to increased annotator experience or any of a number of other time-varying con-
19 founders (Shadish et al., 2002). In another approach, researchers may split annotators
20 into two groups and ask each group to use a different annotation method. Performance
21 statistics are then compared across groups. However, in this design, the causal impact
22 of the annotation design cannot be differentiated from differences in performance due to
23 annotator characteristics. This is particularly problematic when there is a small number
24 of annotators and/or they have not been randomly assigned to their annotation condition
25 (Shadish et al., 2002).

26 In this paper, we demonstrate that the single-case study design can be a key method
27 for building and improving the science of annotation for social scientists and computa-
28 tional linguists alike. This design addresses both time-varying and participant-varying
29 sources of confounding variables by switching the annotation procedure multiple times

1 and comparing outcomes within (rather than across) participants. If the annotation pro-
2 cedure is manipulated many times by the researchers, and the changes in performance
3 track this pattern of manipulation, the researcher can conclude a causal relationship. A
4 key strength of the single case study design is that it can be used with as few as one
5 annotator (Kratochwill et al., 2013). The strong causal validity and low participant re-
6 quirements make it well-suited to empirically testing the efficacy of various annotation
7 techniques in projects relying on just a handful of annotators. For this reason, we argue
8 that researchers should use the single-case study design to guide decisions during the
9 pilot phase of an annotation project and share the results of those studies, increasing the
10 body of empirical knowledge in annotation.

11 This article proceeds as follows. First, we provide an overview of the single case
12 study design for those who may not be familiar. Second, we review key decision points in
13 annotation projects, highlighting points where the single-case study can aid in empirical
14 decision making. Third, we illustrate the application of the design through an applied
15 experiment testing two competing approaches to multi-label annotation projects. Fi-
16 nally, we discuss the generalizability of our results and the strengths and weaknesses of
17 the single-case study design for improving annotation science.

18 **2 The Single-Case Study Design**

19 Single-case study designs have a long history in psychology, dating back to the field’s
20 founders (Perone and Hursh, 2013; Skinner, 1938; Watson, 1925). In contrast to the
21 between-subject design, the single-case study relies on within-subject comparisons, where
22 the participants provide their own control data. The researcher assigns different treatment
23 conditions to the same individual at different points in time while consistently measuring
24 the outcome of interest. If the treatment assignment is manipulated many times by the
25 researcher and the changes in outcomes track this pattern of treatment manipulation, the
26 researcher concludes that the treatment caused the changes in outcomes. This conclusion
27 is warranted when it is difficult to hypothesize confounders that would also produce the
28 observed pattern of effects (Kratochwill et al., 2013). Conclusions from a single case

1 study are primarily drawn from visual analysis of graphs (Kratochwill et al., 2010). To
2 provide evidence of a treatment effect, the graph should demonstrate an unlikely change
3 in the pattern of data that correlates with the researcher’s manipulation of the treatment
4 condition (What Works Clearinghouse, 2019). A stylized example of a convincing single
5 case study is provided in Figure 1.

6 According to the What Works Clearinghouse (a governmental organization that rates
7 the rigor of empirical evidence in education), the single case study design is one of only
8 three designs (including the randomized control trial and the regression discontinuity
9 design) which meet high standards for causal evidence (2019). A strong single case
10 study has the following features: 1) the treatment is manipulated by the researcher, not
11 by the study participants or the environment; 2) the outcome variables are measured
12 systematically and consistently over time; and 3) there are at least three switches in
13 conditions. Together, these conditions reduce the likelihood of confounding variables that
14 produce the same pattern of effects as the manipulation in the treatment assignment.

15 The single case study gets its name from the fact that the design can include as
16 few as one participant. This feature makes it attractive for determining the impact of
17 interventions when the participant pool is small. For example, the design is particularly
18 popular in areas of psychology focused on evaluating treatments for rare or low-incidence
19 diagnoses (Carbone et al., 2013). The limited sample size requirements also make it a
20 low-cost yet causally valid design for researchers making decisions at the beginning of a
21 large annotation project. Annotation projects often include only a handful of annotators,
22 making other causally valid designs, like the randomized control trial, infeasible.

23 The key weakness of the single case design is its potentially limited generalizability.
24 While results can provide evidence of a causal effect for a single individual, this effect may
25 or may not generalize beyond that individual. For this reason, researchers are expected
26 to provide a comprehensive description of participants so that readers may consider the
27 extent to which the impact of an intervention is likely to generalize to their population of
28 interest (Kratochwill et al., 2013). Replicating the single case study design with multiple
29 participants can also provide stronger evidence of a generalizable effect. Further, there

1 are scenarios when a researcher is most interested in identifying a causal effect for their
2 own participants, without any need for generalization. This may occur in clinical cases
3 when an individualized treatment must be chosen, or in annotation projects where the
4 researcher wishes to choose the the most efficient method of annotation for their specific
5 set of annotators. In these cases, the single-case study has few disadvantages.

6 [Figure 1 about here.]

7 **3 Key Questions in the Science of Annotation**

8 Annotation is a complex and multi-part process. As a result, researchers are faced with
9 many decisions in designing and implementing a coding scheme. Here, we focus on
10 two key decisions in an annotation project which can be answered empirically: Who
11 should create the annotations? And how should they do it? For the most part, we
12 sidestep the question of *what* should be annotated, as that decision is wholly dependent
13 on the research question at hand. We'll simply note there is broad agreement that 1) the
14 annotated corpus needs to be representative of the population of interest ([Manning and](#)
15 [Schütze, 1999](#)); and 2) that researchers should create a comprehensive codebook which
16 specifies the definitions of codes and provides examples ([Hovy and Lavid, 2010](#)). Because
17 codes need to be theoretically valid and appropriate for the data at hand, creating the
18 codebook is often an iterative process where the researcher moves back and forth between
19 theory and data before finalizing the code definitions ([Auerbach and Silverstein, 2003](#);
20 [Hovy and Lavid, 2010](#)). Helpfully, researchers can find substantial guidance on creating
21 a codebook from the literature on qualitative research. See, for example, [Auerbach and](#)
22 [Silverstein \(2003\)](#), [Chi \(1997\)](#), and [Shaffer and Ruis \(2021\)](#).

23 **3.1 Who should annotate?**

24 One of the first decisions researchers need to make in an annotation project is who
25 should create the annotations. Researchers may produce the annotations themselves,
26 identify content-area experts to produce the annotations, train undergraduate or gradu-

1 ate students (as is common in academic papers), or rely on untrained annotators from
2 crowd-sourcing platforms like Amazon’s Mechanical Turk (Geiger et al., 2020). There is
3 a general understanding that the cost per annotation resulting from crowd-sourcing can
4 be substantially less expensive than the cost per annotation resulting from content-area
5 experts (Snow et al., 2008; Fort, 2016). However, it is also hypothesized that crowd-
6 sourced annotations will be of lower quality. This hypothesis has been, at least partially,
7 substantiated with empirical evidence. In a comparison of annotations created by ex-
8 pert annotators to those created by crowd-sourced workers, Snow et al. found higher
9 agreement among expert annotators than between expert and non-expert annotators.
10 However, they also found that accuracy can be increased to the level of that achieved by
11 experts by aggregating the annotations of multiple non-experts (2008). Importantly, the
12 accuracy costs of relying on non-expert annotators will be very dependent on the task
13 at hand. The above tests, for example, were completed on tasks requiring only general
14 knowledge of the English language. More specialized tasks may result in lower accu-
15 racy among non-experts. Where relevant, researchers may test this in their own data.
16 Thankfully, there are few causal challenges in identifying the effect of one group of anno-
17 tators versus another. This is because when testing the impact of different annotators,
18 the researcher does not need to worry about annotator characteristics confounding the
19 outcomes; differences between annotators are not confounders but instead the treatment
20 of interest. Thus, so long as the researcher holds other features constant (like time and
21 the annotation task), comparisons of outcomes across participant pools is valid.

22 **3.2 How should the corpus be annotated?**

23 After determining the annotators, researchers need to determine how the annotators
24 will produce their annotations. This involves selecting the annotation procedures and
25 the annotation interface. Regarding the annotation procedures, researchers need to make
26 two key decisions. First, if the annotation task involves multiple codes, should annotators
27 annotate one code at a time, or all at once? Second, if the annotation task involves long
28 documents, what amount of context should annotators use to interpret each text segment?

1 Social science annotators are often advised to work through one document at a time and
2 to consider each piece of text within the context of the full document (Shaffer and Ruis,
3 2021). Computational linguists, on the other hand, are commonly advised to break a
4 complex annotation project down into a series of simple micro-tasks: asking annotators
5 to consider one code at a time and to view the text within just a small context window
6 (Sabou et al., 2014; Hovy and Lavid, 2010). We label these two competing approaches
7 to annotation the *complex* and *simple* annotation schemes.

8 Computational linguists commonly argue that the simple annotation scheme can
9 increase efficiency by placing a lower cognitive load on annotators (Hovy and Lavid,
10 2010; Ide and Pustejovsky, 2017; Sabou et al., 2014). Hovy and Lavid argue, for example,
11 that “though [the simple annotation procedure] compromises on sentence context, [it] is
12 both far quicker and far more reliable: annotators need to hold in mind just one set
13 of alternatives, and become astonishingly rapid and accurate” (2010, p. 10). Similarly,
14 the makers of the popular new annotation software, Prodigy, celebrate the software for
15 allowing annotators to “focus on one task at a time” (Explosion AI, 2017). Though there is
16 an accuracy cost to removing an utterance from its context (Samei et al., 2014), doing so
17 also allows researchers to simplify the annotation task, which is hypothesized to increase
18 annotator efficiency and accuracy enough to make up for performance lost due to lack of
19 context. Further, by decomposing a task into simple yes or no questions, it becomes more
20 feasible to rely on untrained annotators through crowd-sourcing for large-scale projects
21 (Sabou et al., 2014). For example, this is the approach of the Decompositional Semantics
22 Initiative, which decomposes complex linguistic concepts into “straightforward questions
23 on binary properties that are easily answered” by untrained native speakers (White et al.,
24 2016, p.1713). Breaking a multi-label task into multiple simple questions also has the
25 added benefit of flexibility: when annotators code for all codes at once, the codebook
26 becomes brittle. Changes to the coding scheme would require re-annotating all utterances.
27 The simple approach allows for codes to be changed or edited without wasting substantial
28 effort (White et al., 2019). However, this simple approach to annotation is rarely taken
29 by social scientists, either in traditional qualitative research or in text classification. In

1 social science projects, annotators commonly consider a document at a time, annotating
2 for every code in the codebook at once, increasing the cognitive load but also increasing
3 the information available to annotators (see, for example, [D'Angelo et al., 2020](#); [Loksa
4 and Ko, 2016](#)). The applied experiment in this paper demonstrates how the single case
5 study design may be used assess the trade-offs between these two perspectives while
6 controlling for participant and time-varying confounders.

7 After specifying the annotation procedures, researchers need to identify the annota-
8 tion interface, i.e., the software with which the annotators will interact. [Ide and Puste-
9 jovsky](#) identify many potential interfaces including asking annotators to maintain a simple
10 comma-separated-value file, contribute to a SQL database, or use a software specifically
11 designed for annotation ([2017](#)). [Neves and Ševa](#) also provide an extensive review of anno-
12 tation software based on technical criteria (including the cost and easiness of installation),
13 data criteria (including the input and output format of documents), and functionality
14 criteria (including whether the software supports multi-label annotations and document-
15 level annotations). Following these criteria, they recommend three programs that likely
16 meet the needs of most users: WebAnno, brat, and FLAT. Unfortunately, however, there
17 is currently little causally-valid evidence comparing the accuracy and efficiency of annota-
18 tions resulting from competing interfaces. Helpfully, the single-case study design provides
19 a key opportunity to affordably obtain such information.

20 4 Applied Experiment

21 In this section, we demonstrate how the single case study design may be used to inform
22 the development of annotation projects and to answer key questions in annotation science.
23 Specifically, we empirically assess two competing approaches to human annotation. In the
24 simple approach, the annotation task is broken down into short and simple micro-tasks:
25 annotators view short text segments while considering one coding category at a time. In
26 the complex approach, annotators consider all codes at once and consecutively annotate
27 text segments within a full document.

28 The study is situated within a broader educational research project focused on the

1 efficacy of one-on-one coaching for improving teacher practice. The goal of the research is
2 to use text classification to automatically monitor the strategies employed by coaches in
3 their conversations with teachers and teachers-in-training. To this end, a coaching expert
4 developed a coding scheme and codebook by iteratively drawing on coaching research,
5 practitioner resources, their professional experience receiving and providing coaching, and
6 a random sample of coaching transcripts. The initial coding scheme included over 30 po-
7 tential strategies. For the purposes of text classification, we will initially focus on eight of
8 the most common strategies. These include: positive evaluation, observation, suggestion,
9 instruction, demonstration, anticipation, practice, and encouragement. A description of
10 these strategies, along with examples, are provided in Table 1. In a single turn, a coach
11 can employ as many as eight strategies or as few as zero. This means our project involves
12 a multi-label classification task in which there are multiple categories (distinguishing it
13 from a binary classification task) and many can apply at once (distinguishing it from a
14 multi-class task).

15 [Table 1 about here.]

16 4.1 Study Corpus and Participants

17 Our corpus of coaching conversations come from prior studies of the impact of a short
18 (5-minute) coaching intervention on teachers-in-training. For more details on the coach-
19 ing intervention and its effects, see [Cohen et al. \(2020\)](#). All coaching conversations were
20 recorded, professionally transcribed, and segmented by turns of talk. Here, we randomly
21 selected 30 coaching transcripts for piloting annotation, 508 utterances in total. Then,
22 we developed a gold-standard corpus; two coaching experts read the randomly selected
23 transcripts and carefully labelled each coach utterance with the appropriate codes (agree-
24 ment = 0.96, Krippendorff’s alpha = 0.82). Because accuracy was the only priority in
25 the creation of the gold standard corpus, the experts viewed each utterance within the
26 context of the full transcript and took no steps to increase their own efficiency.

27 Four annotators were recruited through the university’s centralized system for hiring
28 undergraduate workers. The job was advertised to students across all schools and ma-

1 jobs at the university. Applicants submitted a resume and short cover letter explaining
2 their interest in the project and participated in a short video interview. While all four
3 annotators had research experience and were in their third or fourth year of study, only
4 two had prior teaching experience or a major within the school of education. Three out
5 of four annotators had prior experience with qualitative coding, specifically. This hiring
6 and recruitment process followed the typical approach in social science research projects
7 (Crittenden and Hill, 1971).

8 In a follow-up experiment exploring mechanisms (discussed later in detail), we sam-
9 pled an additional 20 transcripts, 360 utterances in total. This follow-up experiment
10 was conducted with three of the four annotators. (One annotator graduated from the
11 university and could not participate in the follow-up experiment.)

12 4.2 Annotation Procedures and Interface

13 In line with our annotators' prior technological experiences, we chose a simple annota-
14 tion interface implemented in Excel. Utterances were displayed in one column of the
15 interface and the annotators would enter their codes in a separate column (or columns).
16 Specific annotation instructions depended on whether the annotators were coding under
17 the complex or simple annotation scheme.

18 Under the complex annotation procedure, annotators were asked to code one tran-
19 script at a time and to consider all coaching strategies at once. To this end, their coding
20 interface included one file per transcript. In each file, transcripts were formatted so
21 that each row was a turn-of-talk. Turns-of-talk were kept in the order in which they
22 were spoken, including both coach speech and teacher-in-training speech. For each coach
23 utterance, annotators would select codes from a drop down menu containing the eight
24 coaching strategies and an option for "None of the above." When appropriate, annota-
25 tors could select multiple codes. When annotators finished coding a transcript, they
26 would open the next file to continue coding the next transcript. For an example of this
27 annotation interface, see Figure 2.

28 In the simple annotation scheme, annotators were asked to code for one coaching

1 strategy at a time. Thus, annotators were provided with one file per code (rather than
2 one file per transcript). Again, each row was a turn-of-talk. However, turns-of-talk were
3 presented in random order so that utterances were viewed with only the preceding teacher
4 turn of talk as context. Annotators were then asked to enter a zero or one indicating
5 whether the coach's speech was an exemplar of the target code. Once annotators finished
6 coding all utterances for one coaching strategy, they would open the next file and code
7 the same utterances for the next code. For an example of this annotation interface, see
8 [Figure 3](#).

9 All utterances were coded four times, with two annotators using the complex anno-
10 tation scheme and two annotators using the simple annotation scheme.

11 [Figure 2 about here.]

12 [Figure 3 about here.]

13 **4.3 Measures**

14 For each annotation procedure, we developed analogous methods for measuring efficiency
15 and validity. To assess annotator efficiency, annotators were asked to record their start
16 and end time for each coding file (either the time it took to annotate a transcript or
17 the time it took to annotate all potential exemplars of a coaching strategy). We then
18 converted these values into a measure of time spent per utterance-code, which served as
19 our efficiency metric. In the complex annotation scheme, this was simply the average time
20 it took an annotator to consider the appropriate codes for a coach turn-of-talk. In the
21 simple annotation scheme, this was the summation of the average time it took annotators
22 to consider an utterance for each of the eight coding tasks. Because the simple scheme
23 requires coding the same utterance multiple times (here, eight times), a full picture of
24 coding efficiency requires us to calculate total time spent coding per utterance.

25 To assess validity, we measured the accuracy, precision, and recall of the resulting
26 annotations under each procedure. Accuracy here is defined as annotator agreement
27 with the gold-standard corpus. We measured accuracy by calculating the percent of

1 correctly classified utterance-code pairs; because the annotators classified each turn-of-
2 talk as representative – or not – of eight separate codes, it was possible for an annotator
3 to accurately classify an utterance for one code, but incorrectly classify the utterance for
4 a second code. Because our transcripts were imbalanced (all codes are present in less
5 than 50% of utterances, and some are present in less than 10%), it is also important to
6 measure precision and recall. We measure precision by calculating the proportion of true
7 positive utterance-code pairs out of all positively coded utterances and measure recall
8 by calculating the proportion of true positive utterance-code pairs that the annotator
9 identified as such.

10 **4.4 Study Design**

11 We first randomly assigned four annotators to their starting condition (either the simple or
12 complex annotation procedure). After the first week of coding, annotators were instructed
13 to switch their method of annotation (from the simple to the complex, or vice versa) at
14 the beginning of each of the successive three weeks of coding (see Table 2). In single case
15 study terms, this design is referred to as the ABAB design. It is the switching mechanisms
16 that provides the study with high causal validity; if the impact of switching conditions
17 is strong enough, then it is very difficult to hypothesize alternative explanations for the
18 observed changes in outcomes. Thus, if the simple annotation scheme greatly increases (or
19 decreases) annotation accuracy or efficiency, these changes can be causally attributed to
20 the annotation condition. Our study design meets all of the What Works Clearinghouse
21 standards for a causally-valid single case study ([What Works Clearinghouse, 2019](#)).

22 [Table 2 about here.]

23 **4.5 Statistical Analysis**

24 In this study, efficiency is calculated each week for four weeks and four annotators, result-
25 ing 16 data points in total, an insufficient number of observations for statistical tests of
26 significance, particularly given the dependency structure. However, we have over a thou-
27 sand annotations for each participant: enough to determine whether, for each participant,

1 there is a statistically significant difference in precision, recall, or accuracy depending on
 2 the annotation condition. To this end, we estimate

$$\begin{aligned}
 Y_{ijk} = & \beta_1 \text{Simple}_{ik} \text{Annotator}1_{ik} + \beta_2 \text{Simple}_{ik} \text{Annotator}2_{ik} + \\
 & \beta_3 \text{Simple}_{ik} \text{Annotator}3_{ik} + \beta_4 \text{Simple}_{ik} \text{Annotator}4_{ik} + \\
 & \text{Week}_i \beta + \text{Annotator}_k \beta + \text{Code}_j \beta + \epsilon_{ijk},
 \end{aligned} \tag{1}$$

4 where Y_{ijk} is a binary variable for whether a given turn-of-talk, i , was accurately coded
 5 for code j , by Annotator k . *Annotator* is a vector of indicators for each of the four
 6 annotators, *Week* a vector of indicators for each of the four weeks, and *Code* a vector of
 7 indicators for each of the eight codes in the coding scheme. The coefficients of interest
 8 here are β_1 through β_4 , the average impact of the simple coding procedure for each of the
 9 four annotators.

10 Statistical analyses can also be helpful for summarizing results across participants,
 11 though readers should be careful not to misinterpret the tests of significance. Statistical
 12 inference here is used to make inferences from a sample of utterances to a population of
 13 utterances, not from a sample of annotators to a population of annotators. We summarize
 14 the results across our four participants

$$Y_{ijk} = \beta_1 \text{Simple}_{ik} + \beta_2 \text{Week}_i + \text{Annotator}_k \beta + \text{Code}_j \beta + \epsilon_{ijk}, \tag{2}$$

16 where Y_{ijk} is a binary variable for whether a given turn-of-talk, i , was accurately
 17 coded for code j , for annotator k . *Annotator* is a vector of indicators for each of the four
 18 annotators, *Week* a continuous variable for the week, and *Code* a vector of indicators
 19 each of the eight codes. The coefficient of interest here is β_1 , the average impact of the
 20 simple coding scheme across all four annotators and eight codes.

21 All models were estimated using the `statsmodels` (Seabold and Perktold, 2010) and
 22 `pandas` (Wes McKinney, 2010) packages in Python 3.10.0 (Van Rossum and Drake,
 23 2009). Figures were produced using `matplotlib` (Hunter, 2007) and `seaborn` (Waskom,
 24 2021).

1 4.6 Results

2 The effect of the simple annotation procedure on annotator efficiency is presented in
3 Figure 4. The figure demonstrates that the simple coding procedure took roughly twice
4 as long as the complex coding procedure to produce the same number of codes. In the
5 complex procedure, annotating an utterance for all of the eight codes at once took 35.5
6 seconds on average. Annotating an utterance for a single code took only 8.5 seconds,
7 but this approach requires annotators to read each utterance eight separate times, thus,
8 requiring 68 seconds per utterance to produce the same number of codes as the simple
9 procedure (the sum of the average time spent on each of the eight individual codes). In
10 other words, though reviewing an utterance for a single code took annotators less time
11 than reviewing the utterance for multiple codes, the time spent was not reduced by a
12 factor of eight, which would be required to make the simple annotation more efficient
13 than the complex procedure in this case.

14 From a single case study point of view, Figure 4 provide convincing evidence of
15 causality; manipulation of the treatment condition here is associated with a consistent
16 change in the dependent variable. The effect is visually obvious at each switch in the
17 treatment conditions and is replicated for participant in the study. We see that each
18 individual takes more time when coding under the simple annotation procedure than
19 when coding under the complex annotation procedure. For one individual, this effect
20 seems to be small (Annotator 2), while for the others it is much larger. Crucially, it is
21 very difficult to provide any alternative explanation for the change in times given that the
22 effect is demonstrated at every switch in treatment condition and for every annotator.
23 No other confounding variable is likely to display this same pattern of effects.

24 [Figure 4 about here.]

25 Unlike Figure 4, Figure 5 does not demonstrate a strong consistent impact of the
26 simple annotation procedure on accuracy. While the simple annotation procedure causes
27 a decrease in accuracy for one annotator (Annotator 2), the effect is not convincingly
28 replicated with the other annotators. In Table 3, we summarize the average accuracy for

1 each annotator under the two annotation schemes using Equation 1. While annotator ac-
2 curacy was high across the board (over 95% for each annotator), there are no statistically
3 significant differences in accuracy by annotation procedure for any of the four annotators.
4 When we aggregate these results across annotators using Equation 2, the overall impact
5 of the simple annotation scheme is a small, but significant, decrease in accuracy by half
6 a percentage point.

7 [Table 3 about here.]

8 [Figure 5 about here.]

9 Given the imbalanced nature of our data set, an analysis of precision and recall is
10 important for understanding the impact of the simple annotation procedure. Figure 6
11 demonstrates that the simple annotation procedure caused a decrease in precision for
12 three out of four annotators: on average, a statistically significant negative effect of 3.5
13 percentage points across all four annotators. On the other hand, Figure 7 demonstrates
14 a heterogeneous impact of the simple annotation scheme on recall. While two annotators
15 experienced substantial and consistent impacts of the simple annotation procedure, these
16 effects are in opposite directions (-12 percentage points for Annotator 2 and +10 percent-
17 age points for Annotator 3; see Table 3). These two effects counterbalance one-another,
18 resulting in a very small, non-significant, effect for recall overall. Taken together, the
19 simple annotation procedure increases the time spent coding and reduces precision, but
20 only has a minimal negative impact on overall accuracy.

21 [Figure 6 about here.]

22 [Figure 7 about here.]

23 4.7 Follow-up Experiment to Determine Mechanisms

24 Compared to the complex annotation procedure, the simple annotation procedure is
25 different in two keys ways: 1) the simple procedure provides annotators with less context
26 surrounding an utterance; and 2) the simple procedure asks annotators to review an

1 utterance for a single code at a time. The previous results demonstrated that the simple
2 annotation procedure is less efficient and results in annotations with a lower rate of
3 precision. A natural follow-up question is whether this decrease in efficiency and precision
4 is due to the increase in context or to considering a single code at a time. To test this,
5 we conducted a short follow-up experiment which only varies the context provided to
6 annotators. In both conditions, annotators code for all labels at once (as in the complex
7 procedures), but in one condition ("in-context"), annotators view the utterance within the
8 context of the full transcript. In the other condition ("out-of-context"), annotators only
9 view the preceding utterance. As in the previous design, annotators switched conditions
10 each week.

11 [Table 4 about here.]

12 Figure 8 demonstrates that there is no substantial or consistent efficiency difference
13 for either condition. Thus, context was mostly irrelevant in determining the amount of
14 time coders take to produce annotation. There is also no consistent impact for accuracy
15 or precision. However, when coding in-context, annotators produce annotations with
16 higher recall (by four percentage points; see Figure 9 and Table 4). Taken together with
17 the results of the prior experiment, the results suggest that the amount of context was
18 unimportant for determining annotator efficiency and that the increase in efficiency of
19 the complex annotation procedure was due to annotators considering all codes at once.
20 Interpreting the results of the two experiments in conjunction is a little more complex
21 when considering accuracy, precision, and recall. While the simple annotation procedure
22 reduced precision, a lack of context reduced recall. We discuss potential explanations for
23 these results below.

24 [Figure 8 about here.]

25 [Figure 9 about here.]

5 Interpreting the Results of the Applied Experiments

The above applied experiments tested two key questions in the design of a multi-label annotation task with long documents: should annotators annotate one code at a time, or all at once? And, what amount of context should annotators use to interpret each text segment? Given the results above, we have determined that the best procedure for the annotators in this study is to annotate for all codes at once within the context of a full document (i.e., the complex annotation procedure). We did not find any benefits to accuracy or efficiency resulting from the simple annotation procedure. In total, it took annotators twice as long to code the same data using the simple annotation procedure than using the complex annotation procedure. Whatever cognitive speed was gained by requiring annotators to only consider one code at a time was not enough to outweigh the time it takes to consider the same utterance multiple times.

When considering precision and recall, our results are a bit more complex. While in the main experiment, the simple annotation procedure reduced precision, in the follow-up experiment, a lack of context reduced recall. How can we explain the seemingly distinct results? Consider that recall is a measure of the relationship between the number of false negatives and true positives. If annotators identify fewer true positive, recall will be reduced. It is unsurprising then that reducing context decreases recall because, without context, annotators may not be able to identify every relevant application of a code. On the other hand, annotating for one code at a time likely increases the number of true positives because coders are forced to consider the applicability of every code. Thus, these two opposing mechanisms cancel one-another out in the complex procedure. The remaining decreasing in accuracy in the simple procedure, then, is due to reduced precision: by forcing annotators to consider each code, they are nudged towards (falsely) believing a code is applicable.

Of course, readers should be thoughtful in their consideration of whether the findings of this study generalize to their own context. In particular, there are two dimensions along which generalizability should be considered. First, our study was conducted with undergraduate research assistants, three of whom had prior experience with annotation in

1 other qualitative studies across the university. We might hypothesize that reducing cog-
2 nitive load is more important when annotators lack experience or knowledge of the study
3 context. While our annotators were not content area experts, they were also not novices
4 to the same degree as annotators hired through crowd-sourcing platforms like Amazon
5 Mechanical Turk. However, we believe this project has optimistic implications for the use
6 of crowd-sourcing for social science text classification projects. Crowd-sourcing platforms
7 *necessitate* short simple annotation tasks. MTurkers, for example, expect each task to
8 take a matter of seconds (Sabou et al., 2014). Previous research has demonstrated that
9 crowd-sourced annotators can compete with the accuracy of more traditional annotators
10 (Snow et al., 2008), however, this research does not address the potential loss of accuracy
11 that comes with altering the annotation task so that it may be crowd-sourced. This study
12 demonstrates that while simplifying an annotation task and taking excerpts outside of
13 their larger context may reduce accuracy slightly, it is not to such a degree that social
14 science researchers need to dismiss crowd-sourcing as a possibility. Second, the relative
15 trade-offs of the simple and complex annotation procedures are likely to depend on the
16 coding scheme itself. In particular, the benefits of the simple annotation scheme is likely
17 to vary by the number of codes in the codebook, though the function of this relation-
18 ship is unclear. Further, the amount of text context required for sufficient accuracy will
19 depend on the constructs in the coding scheme. Codes which depend on information
20 provided earlier in conversation will necessitate large context windows. We suggest that
21 in cases where any of these conditions are meaningfully different from the current study,
22 that researchers conduct their own tests. A key strength of the single-case study design
23 is that such tests can be completed quickly and a relatively low cost.

24 6 Conclusion

25 This study demonstrates a straight-forward and low-cost method of testing hypotheses re-
26 garding the design of annotation projects: the single case study design. Given the limited
27 number of participants required to make causal inferences, the single case study design
28 is well-suited to answer annotation questions when projects have only a few annotators.

1 While the randomized control trial would be preferable in the case where an annotation
2 project includes many annotators (say, close to 30), in our experience, researchers rarely
3 hire that many annotators outside the context of crowd-sourcing. The single case study
4 design, on the other hand, is valid with as few as one annotator. Thus, researchers can
5 pilot annotation procedures quickly and cheaply, while also obtaining findings with high
6 causal validity. Though each single case study may only generalize to a subset of annota-
7 tion projects, the relatively low cost of the design means that replicating findings across
8 various contexts is feasible. Thus, we encourage researchers to use the single case study
9 both to inform their own annotation projects and to iteratively improve the evidence base
10 regarding best practices in human annotation.

11 In the past, many text classification papers have neglected to give human annotations
12 the consideration they are due ([Geiger et al., 2020](#)). Indeed, human-annotated training
13 data have been given so little attention that researchers have deemed human-annotated
14 corpora, the “hidden pillars of the domain” ([Fort, 2016](#), p. 9). Thankfully, there is a
15 growing literature on annotation, as evidenced by key texts like [Hovy and Lavid \(2010\)](#)
16 and [Pustejovsky and Stubbs \(2012\)](#), but there are still few empirical tests of the im-
17 pact of annotation conditions on annotation efficiency and quality. We argue, therefore,
18 that researchers should respond to calls for increased attention to annotation quality by
19 incorporating causal evidence into decision-making when designing annotation projects.
20 If human annotations are the “hidden pillars” of text classification, we believe that we
21 can increase the strength and visibility of these pillars through an increased focused on
22 empirical, causally-valid, decision making in annotation.

23 **Supplementary Material**

24 The results of this paper may be reproduced using the scripts and data files uploaded to
25 [Harvard Dataverse](#).

References

- Alyuz N, Aslan S, D’Mello SK, Nachman L, Esme AA (2021). Annotating Student Engagement Across Grades 1–12: Associations with Demographics and Expressivity. In: *Lecture Notes in Computer Science*, 42–51. Springer International Publishing.
- Auerbach C, Silverstein LB (2003). *Qualitative Data: An Introduction to Coding and Analysis*, volume 21. NYU press.
- Bender EM, Friedman B (2018). Data statements for natural language processing: Toward mitigating system bias and enabling better science. *Transactions of the Association for Computational Linguistics*, 6: 587–604.
- Carbone VJ, O’Brien L, Sweeney-Kerwin EJ, Albert KM (2013). Teaching Eye Contact to Children with Autism: A Conceptual Analysis and Single Case Study. *Education and Treatment of Children*, 36(2): 139–159.
- Chi MT (1997). Quantifying Qualitative Analyses of Verbal Data: A Practical Guide. *Journal of the Learning Sciences*, 6(3): 271–315.
- Cohen J, Wong V, Krishnamachari A, Berlin R (2020). Teacher Coaching in a Simulated Environment. *Educational Evaluation and Policy Analysis*, 42(2): 208–231.
- Creswell JW, Miller DL (2000). Determining Validity in Qualitative Inquiry. *Theory Into Practice*, 39(3): 124–130.
- Crittenden KS, Hill RJ (1971). Coding Reliability and Validity of Interview Data. *American Sociological Review*, 36(6): 1073.
- D’Angelo ALD, Ruis AR, Collier W, Shaffer DW, Pugh CM (2020). Evaluating How Residents Talk and What it Means for Surgical Performance in the Simulation Lab. *The American Journal of Surgery*, 220(1): 37–43.
- Dipper S, Götze M, Stede M (2004). Simple Annotation Tools for Complex Annotation Tasks: An Evaluation. In: *Proceedings of the LREC Workshop on XML-based Richly Annotated Corpora*, 54–62.
- D’Mello S (2016). On the Influence of an Iterative Affect Annotation Approach on Inter-Observer and Self-Observer Reliability. *IEEE Transactions on Affective Computing*, 7(2): 136–149.

- 1 Explosion AI (2017). Prodigy: A new tool for radically efficient machine teaching.
2 <https://explosion.ai/blog/prodigy-annotation-tool-active-learning>.
- 3 Fort K (2016). *Collaborative Annotation for Reliable Natural Language Processing: Tech-*
4 *nic and Sociological Aspects*. John Wiley & Sons.
- 5 Gebru T, Morgenstern J, Vecchione B, Vaughan JW, Wallach H, Iii HD, et al. (2021).
6 Datasheets for datasets. *Communications of the ACM*, 64(12): 86–92.
- 7 Geiger RS, Yu K, Yang Y, Dai M, Qiu J, Tang R, et al. (2020). Garbage in, Garbage out?
8 In: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*.
9 ACM.
- 10 Hovy E, Lavid J (2010). Towards a ‘Science’ of Corpus Annotation: A New Method-
11 ological Challenge for Corpus Linguistics. *International Journal of Translation*, 22(1):
12 25.
- 13 Hunter JD (2007). Matplotlib: A 2d graphics environment. *Computing in Science &*
14 *Engineering*, 9(3): 90–95.
- 15 Ide N, Pustejovsky J (2017). *Handbook of Linguistic Annotation*, volume 1. Springer.
- 16 Kratochwill TR, Hitchcock J, Horner RH, Levin JR, Odom SL, Rindskopf DM, et al.
17 (2010). Single-Case Designs Technical Documentation. *Technical report*, What Works
18 Clearinghouse.
- 19 Kratochwill TR, Hitchcock JH, Horner RH, Levin JR, Odom SL, Rindskopf DM, et al.
20 (2013). Single-Case Intervention Research Design Standards. *Remedial and Special*
21 *Education*, 34(1): 26–38.
- 22 Lingren T, Deleger L, Molnar K, Zhai H, Meinzen-Derr J, Kaiser M, et al. (2014). Evalu-
23 ating the Impact of Pre-Annotation on Annotation Speed and Potential Bias: Natural
24 Language Processing Gold Standard Development for Clinical Named Entity Recogn-
25 ition in Clinical Trial Announcements. *Journal of the American Medical Informatics*
26 *Association*, 21(3): 406–413.
- 27 Loksa D, Ko AJ (2016). The Role of Self-Regulation in Programming Problem Solving
28 Process and Success. In: *Proceedings of the 2016 ACM Conference on International*
29 *Computing Education Research*. ACM.

- 1 Manning CD, Schütze H (1999). *Foundations of Statistical Natural Language Processing*.
2 MIT Press, Cambridge.
- 3 Neves M, Leser U (2014). A Survey on Annotation Tools for the Biomedical Literature.
4 *Briefings in Bioinformatics*, 15(2): 327–340.
- 5 Neves M, Ševa J (2021). An Extensive Review of Tools for Manual Annotation of Docu-
6 ments. *Briefings in Bioinformatics*, 22(1): 146–163.
- 7 Perone M, Hursh DE (2013). Single-case experimental designs.
- 8 Pustejovsky J, Stubbs A (2012). *Natural Language Annotation for Machine Learning: A*
9 *Guide to Corpus-Building for Applications*. O’Reilly Media, Inc.
- 10 Sabou M, Bontcheva K, Derczynski L, Scharl A (2014). Corpus Annotation through
11 Crowdsourcing: Towards Best Practice Guidelines. In: *In: Proceedings of the Ninth*
12 *International Conference on Language Resources and Evaluation (LREC’14)*, European
13 *Language Resources Association (ELRA) (2014) 859–866*.
- 14 Samei B, Olney AM, Kelly S, Nystrand M, D’Mello S, Blanchard N, et al. (2014). Domain
15 Independent Assessment of Dialogic Properties of Classroom Discourse. In: *Proceedings*
16 *of the 7th International Conference on Educational Data Mining*, 4.
- 17 Seabold S, Perktold J (2010). statsmodels: Econometric and statistical modeling with
18 python. In: *9th Python in Science Conference*.
- 19 Shadish WR, Cook TD, Campbell DT (2002). *Experimental and Quasi-Experimental*
20 *Designs for Generalized Causal Inference*. Houghton Mifflin.
- 21 Shaffer DW, Ruis AR (2021). How We Code. In: *ICQE 2020* (S Lee, AR Ruis, eds.).
22 Springer, Malibu, CA.
- 23 Skinner BF (1938). *The behavior of organisms: An experimental analysis*. BF Skinner
24 Foundation.
- 25 Snow R, O’Connor B, Jurafsky D, Ng A (2008). Cheap and Fast – But is it Good?
26 Evaluating Non-Expert Annotations for Natural Language Tasks. In: *Proceedings of*
27 *the 2008 Conference on Empirical Methods in Natural Language Processing*, 254–263.
28 Association for Computational Linguistics, Honolulu, Hawaii.
- 29 Van Rossum G, Drake FL (2009). *Python 3 Reference Manual*. CreateSpace, Scotts

- 1 Valley, CA.
- 2 Waskom ML (2021). seaborn: statistical data visualization. *Journal of Open Source*
3 *Software*, 6(60): 3021.
- 4 Watson JB (1925). Experimental studies on the growth of the emotions. *The Pedagogical*
5 *seminary and journal of genetic psychology*, 32(2): 328–348.
- 6 Wes McKinney (2010). Data Structures for Statistical Computing in Python. In: *Proceed-*
7 *ings of the 9th Python in Science Conference* (Stéfan van der Walt, Jarrod Millman,
8 eds.), 56 – 61.
- 9 What Works Clearinghouse (2019). What Works Clearinghouse Standards Handbook:
10 Version 4. *U.S. Department of Education’s Institute of Education Sciences (IES)*, 1–
11 17.
- 12 White AS, Reisinger D, Sakaguchi K, Vieira T, Zhang S, Rudinger R, et al. (2016).
13 Universal decompositional semantics on universal dependencies. In: *Proceedings of the*
14 *2016 Conference on Empirical Methods in Natural Language Processing*, 1713–1723.
- 15 White AS, Stengel-Eskin E, Vashishtha S, Govindarajan V, Reisinger DA, Vieira T, et al.
16 (2019). The universal decompositional semantics dataset and decomp toolkit. *arXiv*
17 *preprint arXiv:1909.13851*.

1 **List of Figures**

2 1 Stylized example of ABAB single case study design with one participant
3 and clear causal impact. Outcomes may either be single observations taken
4 from the participant or average outcomes across many observations of the
5 same participant. The strongest single case studies are also replicated
6 multiple times with more than one participant. 26

7 2 Complex annotation procedure interface. All coach and teacher utterances
8 in a given transcript were included in the order in which they were spoken.
9 Annotators considered one transcript at a time. 27

10 3 Simple annotation procedure interface. Coach utterances were presented in
11 randomized order along with the preceding teacher utterance. Annotators
12 considered one code at a time. 28

13 4 Average total annotation time per utterance as a function of the annotation
14 procedure. The y-axis of the figure displays the total average annotation
15 time per utterance. Under the complex annotation scheme, this is the aver-
16 age time it takes the annotators to consider the relevance of the eight codes
17 all at once for a given utterance. Under the simple annotation scheme, this
18 is the average total time it takes annotators to consider the relevance of
19 the eight individual codes one at a time. 29

20 5 Accuracy as a function of the simple versus complex annotation procedure.
21 The y-axis of the figure displays average accuracy across all utterance-code
22 pairs. 30

23 6 Precision as a function of the simple versus complex annotation procedure.
24 The y-axis of the figure displays average precision across all utterance-code
25 pairs. 31

26 7 Recall as a function of the simple versus complex annotation procedure.
27 The y-axis of the figure displays average recall across all utterance-code
28 pairs. 32

29 8 Average total annotation time per utterance as a function of the context
30 provided to annotators. In both procedures, annotators code for eight
31 codes at once, but while the "In Context" procedure shows all utterances
32 in order, the "Out of Context" procedure present coach utterances in ran-
33 domized order. The y-axis of the figure displays the total average annota-
34 tion time per utterance. 33

35 9 Recall as a function of the context provided to annotators. In both proce-
36 dures, annotators code for eight codes at once, but while the "In Context"
37 procedure shows all utterances in order, the "Out of Context" procedure
38 present coach utterances in randomized order. The y-axis of the figure
39 displays average accuracy across all utterance-code pairs. 34

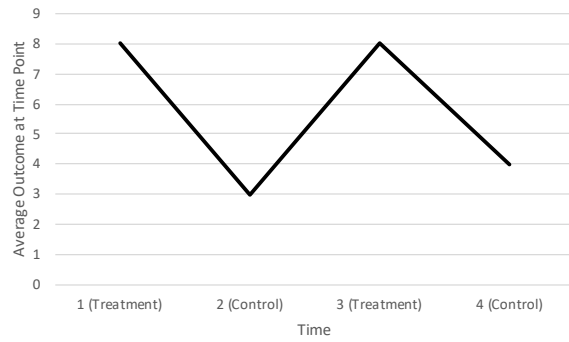


Figure 1: Stylized example of ABAB single case study design with one participant and clear causal impact. Outcomes may either be single observations taken from the participant or average outcomes across many observations of the same participant. The strongest single case studies are also replicated multiple times with more than one participant.

Utterance ID	Speaker	Text	Code 1	Code 2	Code 3	Code 4
426	Coach	So first, how'd you feel?				
427	Teacher	It was good at first, but went downhill from there.				
		Which is totally fine, like, and this is something that just takes practice and if we're going to get a feel for it in your second attempt here. When you say it went downhill, like what do you think?				
428	Coach					
429	Teacher	Ethan was distracting everyone else and I spent most of my time talking to Ethan.				

1 Tellback Positive Evaluation

2 Tellback Observation

3 Tellforward Suggestion

4 Tellforward Instruction

5 Tellforward Demonstration

6 Askforward Anticipation

7 Practice

8 Rapport Encouragement

Teacher

NA

Figure 2: Complex annotation procedure interface. All coach and teacher utterances in a given transcript were included in the order in which they were spoken. Annotators considered one transcript at a time.

STRATEGY: Positive Evaluation			
Utterance ID	Preceding Teacher Text	Coach Text	Code
603	Mm-hmm.	The other thing I do want to point out is there were quite a few times throughout your past stimulation when um you did provide various specific um instructions for an attempt to redirect the misbehavior.	0
416	I think I did a good job of like asking students to explain or like pullback in the text to explain their answer instead of just what they think.	I heard you say, "what in the text make you think that?" or like "could you read me this part of the text?" and that's like a great probe for like textual evidence, that was awesome. Are there things that you, like, think you could have done differently in terms of feedback?	1

Figure 3: Simple annotation procedure interface. Coach utterances were presented in randomized order along with the preceding teacher utterance. Annotators considered one code at a time.

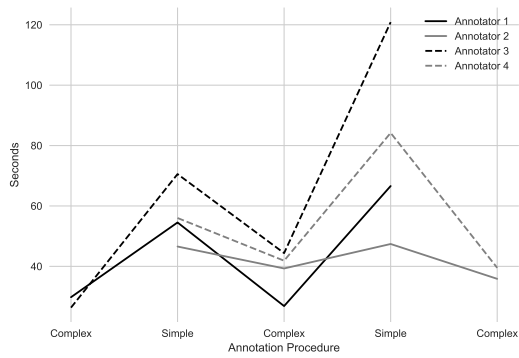


Figure 4: Average total annotation time per utterance as a function of the annotation procedure. The y-axis of the figure displays the total average annotation time per utterance. Under the complex annotation scheme, this is the average time it takes the annotators to consider the relevance of the eight codes all at once for a given utterance. Under the simple annotation scheme, this is the average total time it takes annotators to consider the relevance of the eight individual codes one at a time.

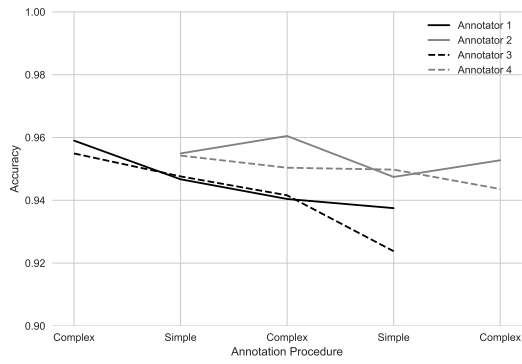


Figure 5: Accuracy as a function of the simple versus complex annotation procedure. The y-axis of the figure displays average accuracy across all utterance-code pairs.

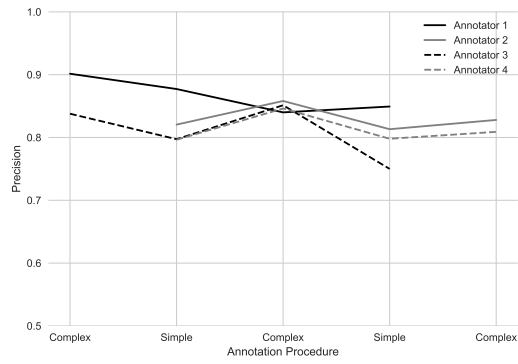


Figure 6: Precision as a function of the simple versus complex annotation procedure. The y-axis of the figure displays average precision across all utterance-code pairs.

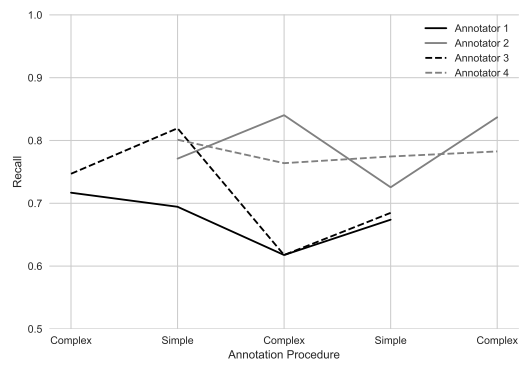


Figure 7: Recall as a function of the simple versus complex annotation procedure. The y-axis of the figure displays average recall across all utterance-code pairs.

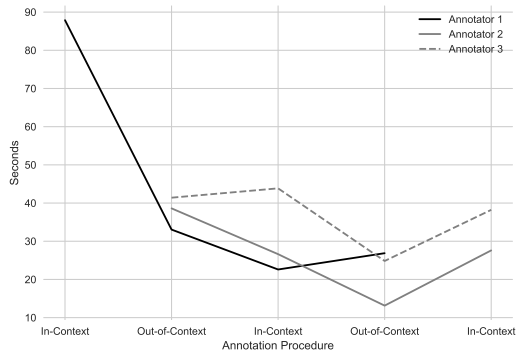


Figure 8: Average total annotation time per utterance as a function of the context provided to annotators. In both procedures, annotators code for eight codes at once, but while the "In Context" procedure shows all utterances in order, the "Out of Context" procedure present coach utterances in randomized order. The y-axis of the figure displays the total average annotation time per utterance.

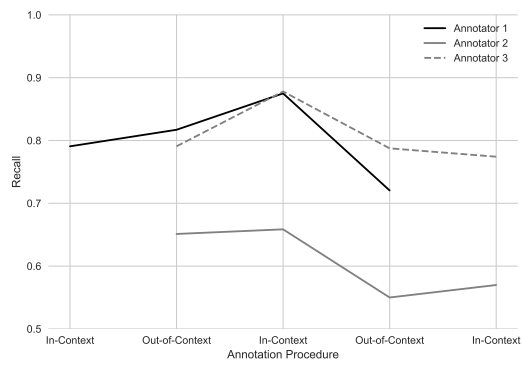


Figure 9: Recall as a function of the context provided to annotators. In both procedures, annotators code for eight codes at once, but while the "In Context" procedure shows all utterances in order, the "Out of Context" procedure present coach utterances in randomized order. The y-axis of the figure displays average accuracy across all utterance-code pairs.

1 List of Tables

2	1	Coding scheme.	36
3	2	Study design and annotation procedure assignments.	37
4	3	Impact of the simple annotation procedure on accuracy, precision, and recall.	38
5	4	Impact of the lack of context (when annotating for all codes at once) on accuracy, precision, and recall.	39
6			

Table 1: Coding scheme.

Strategy	Definition	Example
Positive Evaluation	Positive judgement about a teacher's skills or practice	<i>You are so kind and engaging with him!</i>
Observation	Specific information about the students or teacher based on the coach's observation	<i>You tend to ask kids to raise their hands a lot.</i>
Suggestion	Explicit proposal that the teachers can or should make to their instruction	<i>One thing you could do is to try to avoid a negative tone of voice.</i>
Instruction	Information that helps a teacher understand the importance or purpose of an instructional strategy	<i>Being more specific with your redirections ensures that your students understand your expectations and can follow them.</i>
Demonstration	A specific demonstration of how to implement and instructional strategy	<i>A calm tone of voice would sound like, "Ethan, please be quiet".</i>
Anticipation	A question that prompts the teacher to elaborate on the consequences of an instructional strategy	<i>What do you think would happen if you asked students to raise their hands?</i>
Practice	Dialogue where the coach facilitates a role-play activity	<i>We're going to practice. I'll pretend to be a student and I want you to redirect me.</i>

Table 2: Study design and annotation procedure assignments.

Annotator	Week 1	Week 2	Week 3	Week 4
1	Complex	Simple	Complex	Simple
2	Simple	Complex	Simple	Complex
3	Complex	Simple	Complex	Simple
4	Simple	Complex	Simple	Complex

Table 3: Impact of the simple annotation procedure on accuracy, precision, and recall.

	Accuracy ($M = 0.95$)	Precision ($M = 0.75$)	Recall ($M = 0.83$)
Annotator 1*Simple Scheme	-0.003 (0.008)	-0.003 (0.037)	0.041 (0.041)
Annotator 2*Simple Scheme	-0.011 (0.008)	-0.035 (0.036)	-0.119** (0.036)
Annotator 3*Simple Scheme	-0.006 (0.009)	-0.052 (0.04)	0.103** (0.04)
Annotator 4*Simple Scheme	-0.001 (0.008)	-0.043 (0.038)	-0.014 (0.04)
Average Impact Across Annotators	-0.005* (0.003)	-0.035** (0.012)	0.003 (0.015)

Note. $N = 508$. The first four rows of the table represent the impact of the simple annotation procedure on each of the four annotators' accuracy, precision, and recall, estimated using Equation 1. The final row represents the average impact of the simple annotation procedure across all four annotators, estimated using Equation 2. *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$.

Table 4: Impact of the lack of context (when annotating for all codes at once) on accuracy, precision, and recall.

	Accuracy ($M = 0.95$)	Precision ($M = 0.74$)	Recall ($M = 0.85$)
Annotator 1*No Context Scheme	0.005 (0.009)	-0.029 (0.044)	-0.03 (0.052)
Annotator 2*No Context	0.005 (0.009)	-0.029 (0.044)	-0.03 (0.052)
Annotator 3*No Context	0.005 (0.009)	-0.009 (0.048)	-0.054 (0.047)
Average Impact Across Annotators	-0.002 (0.003)	0.007 (0.016)	-0.043* (0.022)

Note. $N = 360$. The first three rows of the table represent the impact of the lack of context, when coding for eight codes at once, on accuracy, precision, and recall for each of three annotators, estimated using Equation 1. The final row represents the average impact of lack of context across all three annotators, estimated using Equation 2. *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$.