A streamlined approach to online linguistic surveys*

Michael Yoshitaka Erlewine and Hadas Kotek, MIT {mitcho,hkotek}@mit.edu

Abstract

More and more researchers in linguistics use large-scale experiments to test hypotheses about the data they research, in addition to more traditional informant work. In this paper we describe a new set of free, open-source tools that allow linguists to post studies online. These tools allow for the creation of a wide range of linguistic tasks, including linguistic grammaticality surveys, sentence completion tasks, and picture-matching tasks. Our tools further help streamline the design of such experiments and assist in the extraction and analysis of the resulting data. Surveys created using the tools described in this paper can be posted on Amazon.com's Mechanical Turk, a popular crowdsourcing platform which mediates between 'Requesters' who can post surveys online and 'Workers' who complete them. This allows many linguistic surveys to be completed within hours or days and at low costs of \$50 or less. Alternatively, researchers can host these randomized experiments on their own servers using a supplied server-side component.

1. Introduction

In recent years, researchers in linguistics have begun to integrate experimental components into their studies, in addition to more traditional informant work. With crowdsourcing tools such as Amazon's Mechanical Turk (henceforth: AMT), researchers can develop and test hypotheses on a wide variety of naïve speakers within days at very low cost. The purpose of this paper is to describe the process of creating an experiment, uploading it onto AMT and obtaining results using originally-developed tools that allow researchers who have little programming experience to create linguistic surveys. These tools, as well as several example templates corresponding to different kinds of tasks, were written under a liberal, open-source license and are distributed for free online at https://github.com/mitcho/turktools.

The tools described here are designed for use with AMT, a popular online crowdsourcing website where companies or individuals (called *requesters*) can post small tasks (called Human Intelligence Tasks, or *HITs*) that cannot easily be automated, and therefore require human workers (called *workers*) for completion. HITs are generally small in nature (such as identifying the contents of an image), and generally high in quantity (it is not unusual for requesters to post thousands of tasks in a single batch). Requesters generally pay very little per HIT (e.g., several US cents) and retain the ability to Accept or Reject the results of each HIT before AMT sends payment to the worker. HITs can be posted using an online interface (www.mturk.com), and

1

^{*} Acknowledgements to be added

their results are easily analyzable using mainstream statistical software (for more information, see Fort et al. 2011; Sprouse 2011; Gibson et al. 2011).

In the next section, we will briefly discuss the value of experimental methods and, in particular, crowdsourcing for theoretical linguistics. In the remainder of the paper, we describe *turktools*: a set of tools that enable easy use of online surveys for obtaining linguistic data, and a suggested workflow for their use.¹

- (1) **Skeletons**: HTML skeletons for various experiment types, including grammaticality judgment tasks with scales or yes-no choices, picture-matching tasks, tasks involving audio files, forced-choice completions.
- (2) **Templater**: A Python script that creates HTML templates out of skeletons.
- (3) **Lister**: A Python script that takes experimental items, including fillers, and organizes them into randomized, counterbalanced lists for presentation in the online survey.
- (4) **Simulator**: A Python script that takes a template and randomized items in an AMT-style lists file and displays a version of the survey as it would appear in the online survey.
- (5) **Decoder**: A Python script that takes a survey results file from AMT and decodes it into a format that can easily be analyzed by any statistical software.
- (6) **Analysis.r**: A sample R script that assists in preparation and analysis of the online survey results.

While these tools were built with the AMT platform in mind, they by no means require the use of AMT. At the end of the paper we describe a server-side tool that allows for the hosting of AMT-style surveys online on a researcher's own server. For brevity, throughout the paper we describe the process of creating an experiment as if the researcher will upload the survey to AMT, but the majority of the process also applies to surveys hosted on a researcher's own server.

In addition to describing this workflow, we will touch on both practical and philosophical aspects of experimental design and analysis throughout.

2. The value of crowdsourcing in theoretical linguistics

Current research in theoretical linguistics relies heavily on *direct-informant* work – that is, obtaining linguistic data from one or several native-speaker consultants for a given language. This is often done by presenting the consultant with a sentence in a context or paired with a certain desired meaning and asking whether the sentence in appropriate in the given context (Chomsky, 1965; Schütze, 1996). This method has been criticized for involving (a) relatively

¹ At the time of writing, these tools require the use of Python 2.6.x or 2.7.x, available at http://python.org. However, the tools described here and their prerequisites and usage are subject to change. Please consult the latest information at https://github.com/mitcho/turktools before using these tools.

few participants, (b) a relatively small number of target stimuli, and (c) cognitive biases on the part of the researcher and participants (Cowart 1997; Edelman and Christiansen 2003; Featherston 2005; Ferreira 2005; Gibson and Fedorenko 2010; Marantz 2005; Myers 2009a,b; Schütze 1996; Sprouse 2009; Wasow and Arnold 2005).

Linguists have consistently offered strong arguments in response to these criticisms (e.g., Phillips and Lasnik, 2003; Marantz, 2005; Culicover and Jackendoff, 2010) and have reported formal experimental results that corroborate many informal experimental results (Featherston, 2005; Phillips, 2009). Many judgments obtained through informal methods and presented in journals and textbooks have been experimentally replicated using diverse experimental methods and tasks, thus showing that many of the criticisms cited above are unwarranted (Cable and Harris 2011, Munroe et al. 2010, Sprouse and Almeida 2010, 2012, ms., Sprouse 2011, Sprouse et al. 2012).

Following this latter work, we believe that experimental methods are not necessary in order to substantiate every linguistic claim made in the literature. Nonetheless, such methods are often useful and helpful in the process of theory-building. With crowdsourcing tools as AMT, researchers can develop and test hypotheses on a wide variety of naïve speakers within days at very low cost. Therefore, although we do not advocate the use of these methods in all cases and at all costs, we believe that it is vital for theoretical linguists to develop the skills that would allow them to use such methods when necessary.

The tools we describe here are designed specifically for use with AMT. We recognize that some concerns have been raised about the nature of the AMT participant pool in the context of linguistic experiments (see Fort et al. 2011), in particular as to the number of workers who complete most HITs that are posted on AMT. According to Fort et al. (2011), 20% of the workers complete 80% of the tasks, thus casting doubt on the diversity of the linguistic data obtained through AMT surveys. From our own experience with AMT, however, we believe that such concerns are unfounded. In all the AMT experiments conducted at the LAB between November 2010 and April 2013, a total of 4635 unique workers participated in at least one study; 643 (14%) participated in more than one study, 152 (3%) participated in more than two studies, and only 15 (0.3%) participated in more than five. Thus, we believe that the diversity of our data is not jeopardized by the tendencies described in Fort et al. (2011).

3. Formulating a research question and choosing an experimental design

Before describing in more detail the technical tools that allow one to prepare a linguistic survey and upload it on AMT, the first task the linguist must address is formulating a clear hypothesis

-

² Data collected on April 24, 2013. The vast majority of experiments were on English and restricted IP addresses of workers to within the US.

that can be tested using experimental methods. The simplest approach, and the one that we will be demonstrating in this paper, is crossing *factors* that have two *levels* each. Factorial designs have proven useful in linguistics research, and are commonly used in many areas of experimental cognitive science. 2×2 designs, in particular, are common because of their simplicity and the relative ease of the analysis of their results. In this paper we will illustrate the use of the tools presented here through an example with a 2×2 design. We leave it to the readers to scale this design up or down to suit their own needs.

In the remainder of this paper, we illustrate the use of factorial designs with a well-known example from the literature: *there*-constructions have been argued to show a *definiteness effect*, whereby a definite or quantified DP cannot appear in a *there*-construction, (8a-c) (Milsark 1974, 1977 and much subsequent work).

- (7) a. There is a boy in the room
 - b. *There in the boy in the room
 - c. *There is every boy in the room

In what follows we restrict our attention to the contrast between (7a-b). We may formulate the hypothesis that *definite DPs are ungrammatical in a there-construction*. To test this, we will compare sentences that differ in whether they contain a definite or an indefinite DP. We will cross this *factor* with whether or not a sentence contains a *there-*construction, thus yielding a *baseline* for the sentences of interest in our survey. A sample item set for a simple survey testing the grammaticality judgments reported in the literature for the *there-*construction is given in (8), with predicted grammaticality judgments indicated.

(8) a. A boy is in the room.	–def, –there
b. The boy is in the room.	+def, –there
c. There is a boy in the room.	-def, +there
d. *There in the boy in the room.	+def, +there

In many cases, testing multiple similar items, differing only in the *lexicalizations* and otherwise systematically manipulating the factors of interest, may further the validly of the experimental results. *Turktools* straightforwardly supports the generation of multiple items within each experimental condition. We will discuss below how randomization and list creation is supported. We refer the reader to the literature on internal validity and external validity for suggestions on how to correctly define one's research goals, whether multiple test items are adequate for a particular study, how many items to construct, and how to do so (see e.g. Gelman and Hill 2007, Myers 2009a and references therein).

The tools that we provide in this paper support diverse types of experiments, some of which we describe below. For each type of experiment we describe here, we provide a *skeleton* HTML file for creating such an experiment. We will discuss skeletons in more detail in Section 4. For other

types of experiments used in the psycholinguistic literature as well as recommendations for design and data analysis, see Schütze and Sprouse (to appear) and citations therein.

3.1 Acceptability judgment task

Perhaps the most commonly used experimental task using in syntax and semantics research is the acceptability judgment task where participants are asked to judge the acceptability or naturalness of a sentence, or of a sentence-meaning pair. The dependent measure in this task is a rating on some scale, where the ends of the scale may correspond to 'acceptable'-'unacceptable', 'natural'-'unnatural', 'grammatical'-'ungrammatical', 'good'-'bad', etc. Alternatively, a fixed Likert scale, often with five or seven points, may be used. A similar task known as magnitude estimation requires participants to compare the acceptability of target sentences to some predetermined reference sentence (Bard et al. 1996).³

The acceptability judgment task is perhaps the easiest language task to implement for online surveys. We provide skeletons for both binary forced-choice (e.g. 'yes' or 'no') and Likert-scale implementations. Both types of skeletons are easily adaptable to allow for different answer choices, such as 'ves'-'no'-'don't know', for a different number of options on a Likert scale, and for magnitude estimation tasks.

3.2 Forced-choice completion task

Another simple task is one in which a sentence or context is presented to the participant, where the task is either to choose which word or phrase better completes a sentence, or which whole sentence is better suited to describe a given situation. An example word completion task is given below in Figure 1. The dependent measure is the rate of selection of each choice in comparison with its competitor.⁴

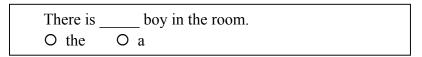


Figure 1. A sample word completion task

3.3 Picture-matching tasks

For tasks in which a visual scene must be presented, a picture-matching task is often used. In this task, participants are shown a picture and are asked to choose which among several options best describe the sentence, or alternatively they are presented with a sentence and several pictures and must choose which among the pictures best matches the description in the sentence. The

³ See Cowart 1997, Keller 2000, Featherston 2005, Sprouse 2011, Weskott and Fanselow, 2011, among others, for more on this task, including its advantages and potential criticisms.

4 Our supplied skeletons support choices introduced with buttons below the sentence, as in Figure 1, or with a drop-

down menu.

dependent measure is the rate of choice of a given sentence or picture compared to its competitors. In another variant of the task, participants are presented with a picture and a sentence that describes it, and must decide whether or not the sentence is a good description of the picture. The picture-matching skeleton we provide also allows for an easy implementation of the 'Covered Box' technique, which is particularly suited for detecting less preferred interpretations of ambiguous sentences (for details, see Huang, Spelke and Snedeker, ms.; Pearson, Khan and Snedeker 2010).

3.4 Additional modifications

For the experiment types described above (and others that users of our tools might develop), several modifications can be easily made.

- An audio file may be added.
- A context may be added before or after the target sentence.
- A comprehension question may be added either before or after the target sentence.
- For *cognitive load* paradigms, the target task may be interrupted by a secondary task.
- Items may be presented one at a time instead of all together on the same page.
- The sequence of events with which a single item is shown can be controlled for example, it is possible to first show a picture, then mask the picture and ask a question about the picture. This setup is also useful for *cognitive load* paradigms.
- Timing information may be logged for each item or parts of an item that have been viewed.⁵

4. Templater: Creating an HTML template for use on Mechanical Turk

Once a research question is formulated and an experiment type has been selected to test it, the presentation of the survey can be determined. Surveys are presented on AMT using an HTML template. We recommend creating a template in two steps: (a) choosing an appropriate *skeleton* out of the ones offered on https://github.com/mitcho/turktools and editing it to fit your study, and then (b) creating a template from the skeleton using the *Templater*.

Skeletons are recipes for different experiment types that abstract away from the number of items that will be presented in a given study. They are stripped-down versions of HTML templates which should be edited to fit each particular design or experiment. We recommend the use of a

_

⁵ This might require custom JavaScript programming in the template, and furthermore such timing information may not be very reliable. Our own templates utilize the jQuery JavaScript library (http://api.jquery.com/), and we recommend its use for such custom programming. See Ibex, "Internet Based Experiments", http://spellout.net/ibexfarm/ for a tool specifically designed for conducting timing-sensitive experiments online.

programming text editor to edit these files.⁶ For each experiment type we surveyed in Section 3, we provide a corresponding skeleton, (9).

(9) Supplied skeleton files:

grammaticality judgment yes-no task binary.skeleton.html binary-image.skeleton.html picture matching yes-no task completion.skeleton.html word-completion task c. image-choice.skeleton.html picture matching task d. grammaticality judgment scale task likert.skeleton.html f sentence-choice.skeleton.html sentence-completion task completion-menu.skeleton.html completion task with drop-down menu

A skeleton is itself an HTML file, but contains *substitution tags* which will be filled in by the Templater. These tags are all wrapped in double curly braces, i.e. {{...}}.

The skeletons that we provide all share the same basic structure, illustrated in Figure 2. The very top of the skeleton has a comment to remind users that a template should be created out of the skeleton to be uploaded on AMT.⁷ Each experiment will be assigned a code, is used to instruct workers to only complete one survey of a certain type.⁸ Next are the instructions for the survey, including any practice items, as well as a consent statement and contact information for the experimenters.⁹ Finally, there is text requesting workers to turn on JavaScript in their browser. This enables a counter that is included in all the skeletons which helps workers make sure that they have answered all the questions in the survey.¹⁰

_

⁶ Good, free options include *Notepad++* for Windows and *TextWrangler* for Mac.

⁷ Comments, which will be completely ignored and removed by the Templater, are tags with an exclamation point, as in $\{\{1,\ldots\}\}$.

⁸ This is often desired in linguistic experiments, since collecting multiple judgments from the same language consultant may result in skewed results of a study, if the researcher mistakenly treats these observations as independent.

⁹ Approval for your proposed experiments from the Institutional Review Board (IRB) at your institution may be required before you run them on AMT. Please consult your institution's IRB.

¹⁰ Researchers can make the use of JavaScript obligatory, for example if JavaScript is used for critical manipulations such as hiding and revealing stimuli or recording timing information. For the template skeletons which we provide, it is possible for participants to take part in surveys even if JavaScript is not enabled in their browser. The result is that they will not see the counter, but other functionality of the survey will be maintained.

```
{{! This is a template skeleton; use templater.py!}}
Survey Code: {{code}}
PLEASE COMPLETE AT MOST ONE {{code}}} SURVEY. YOU WILL NOT BE PAID
FOR COMPLETING MORE THAN ONE SURVEY WITH THIS CODE.
Instructions
This questionnaire presents {{total number}} English sentences. Each sentence contains a gap
and there are two options below the sentence for what should go in that gap. Choose the option
that sounds more natural to you. Here is an example:
               There is ____ boy in the room.
               O the O a
In this example...
Consent Statement: ...
If your browser has JavaScript turned on, a counter will be displayed at the bottom of the page
indicating how many questions have been answered. It is highly recommended that you turn on
JavaScript and use this tool before submitting to ensure that all questions have been
answered and you can receive payment.
{{#items}}
{\text{number}}. {\text{field\_1}} \_ {\text{field\_2}}.
              O {{field 3}} O {{field 4}}
{ {/items } }
```

Demographic questions, questions about native language.

0 questions (out of {{total_number}} total) have been answered. If you submit now, you will not be paid.

After submitting this HIT, do NOT submit another HIT with survey code {{code}}. You will not be paid for completing more than one survey with this code.

Figure 2. An example skeleton file

Notice the substitution tags in this skeleton. The templater will replace {{code}} with the experiment's unique code, and {{total_number}} with the number of items presented in the experiment.

The main body of the skeleton is the items block, beginning with $\{\{\text{#items}\}\}\$ and ending with $\{\{\text{/items}\}\}\$. The items block contains one sample item of the shape that all items in the survey will take. When a template is created out of the skeleton, this block will be duplicated as many times as there are items in your survey. The item contains the tag $\{\{\text{number}\}\}\$, which will be replaced with the item number in the experiment, as well as a number of $\{\{\text{field}_n\}\}\$ tags. In this case, each item is made up of four fields: fields 1 and 2 are the parts of the sentence before and after the gap and fields 3 and 4 correspond to possible answer options. Each of these $\{\{\text{field}_n\}\}\$ tags will be replaced by a different "field" in the items file. The makeup of the items file will be described in the following section.

Below the items block are some demographic questions. We always ask participants to indicate their native language and any other languages they speak, although we clarify that payment is not contingent on their answers to these questions. Researchers may wish to add questions about gender, age, origin, etc. Finally, at the bottom of the skeleton is a counter to help workers ensure that they reply to all questions before submitting their survey. Below the counter is text warning participants not to accept a second HIT with the same code after submitting their survey.

Note that all parts of the skeleton can be edited, moved around or removed to suit the preferences of the researcher and the study. We recommend saving the changes to the skeleton in a dedicated folder that will contain all files pertaining to the same experiment. At least some changes will only have to be made once – for example, the consent statement should be the same in all surveys conducted under the same IRB approval.

After choosing a skeleton and editing it, a template must be created out of the skeleton using the *Templater*. The script *templater.py* should be saved in the same folder as the skeleton. The *templater.py* script and all other Python scripts described in this paper can be run in the following manner:

(10) Executing a Python script (here templater.py):

- a. <u>UNIX shell:</u> move (cd) to the directory that contains the script and execute python templater.py.
- b. Mac OS X: Right-click on the file in the Finder and choose Open With... > Python Launcher.
- c. Windows: Double-click on the file. 11

¹¹ For Mac OS X and Windows, these instructions assume that Python was installed using the standard Python binary packages provided at http://python.org. If not, use the UNIX shell instructions via the Terminal on Mac, and open using the python application on Windows.

The *Templater* will ask for three pieces of information, (11). Once these are provided, a template file will be created. The template file will be located in the same folder as the skeleton file and the *templater.py* script.

- (11) a. the skeleton file name: one of the skeletons provided here or others you create yourself.
 - b. the number of items in your survey: including all experimental and filler items but not practice items, which are coded in the skeleton.
 - c. a survey code: any letter-number combination you choose.

Note that the *Templater* requires an items file that is formatted in a way that matches the skeleton you chose. That is, each item must have at least as many fields as there are fields in the sample item in the skeleton that you have chosen. 12 Thus, the choice of design and skeleton will dictate the format of your items file. We discuss the items file in Section 5.

5. Creating an items file

Once a design and a template have been decided on for your survey, items need to be created. For a simple 2×2 experiment crossing two factors (here: definiteness, there-construction) with two levels each (here: +/- definite, +/- there-construction), each target item set will contain four different *conditions*. A sample target *item set* for a sentence judgment task will look as in (12).

(12) A sample item set:

	Stimulus:	Condition name:	Factor values:
a.	A boy is in the room.	indef-no.there	-definite, -there-condition
b.	The boy is in the room.	def-no.there	+definite, -there-condition
c.	There is a boy in the room.	indef-there	-definite, +there-condition
d.	There is the boy in the room.	def-there	+definite, +there-condition

Turktools uses a naming convention whereby the condition name specifies the chosen value for each of the factors, separated by hyphens (13). Names of this form will be assumed by the decoder, which we will discuss in section 9.13 For our there-construction experiment, this results in conditions names as in the middle column of example (12), where the settings of the two factors are shown in the third column and the first column shows the actual sentence that we would like participants to rate.

(13)The structure of a condition name:

<factor 1 value>-<factor 2 value>-<factor 3 value>-...

¹² It is possible to have more fields in the items file than in the template, but not the other way around. The use of hidden fields will be discussed in Section 5.

13 It is possible to use other naming schemes for condition names. In such a case, however, part of the decoding

process, which is automated in the *decoder*, will have to be done manually by the researchers. See footnote 21.

An experiment is composed of different sections, where normally at least one section will be used as *fillers*, and the others will be *targets*. This designation is made when using the Lister (section 6). Each item in the items file has two parts, each beginning on a new line: 14

- (14) a. **The item header**. The item header consists of four components, separated by spaces.
 - (i) the symbol #;
 - (ii) the name of the section (e.g., target in the example in Figure 3);
 - (iii) the item set number within the *section* (we recommend 1, 2, 3, ...);
 - (iv) the *condition* name (e.g., indef-no.there, def-no.there, indef-there, def-there in our example). The condition names are not constrained in any way, but it is most helpful to use labels that reflect the experimental design, and in particular to describe the setting of the *factors* of interest for that item.
 - b. The item body. The item body consists of fields, each on its own line. 15 Line n corresponds to the text that will be substituted for the text $\{\{field_n\}\}\$ in the skeleton file. A field may specify a sentence to be judged, choices for completions, a picture, an audio file, a context, a comprehension question, etc.

It is possible to add hidden fields to the body of an item by adding extra lines that do not have corresponding substitution tags in the HTML template. This is a way to specify additional information for a given item that is useful for the researcher but should not be available to the participants. For example, if a skeleton refers to {{field_1}} through {{field_4}} but not {{field 5}}, then anything that occurs in the 5th line (and onward) of an item's body will not be displayed on the screen and participants will not have access to this information. However, the data will still be part of the results file created by AMT after the experiment is completed and can be used as part of the data analysis. For example, the correct answers to comprehension questions or expected answers to filler items may be specified in this way.¹⁶

The format for a sample item set in a there-construction study is given in Figure 3. For illustration, the experiment here pairs each sentence with a picture. Each item consists of two

¹⁴ This structure is very similar to that used for items files in Linger, written by Doug Rohde (http://tedlab.mit.edu/ ~dr/Linger/), and Gibson et al.'s (2011) Turkolizer. Note that our terminology differs slightly: Gibson et al refers to a group of trials together forming an item, which corresponds to our notions of individual items which are grouped together into item sets.

Furthermore, the structure of the item body in our Lister format is much more flexible, allowing for an arbitrary number of fields. Unlike in Linger and Gibson et al.'s Turkolizer, comprehension questions are not given a special status but are instead treated like any other field that could be added to an item. If a comprehension question is given or if some filler items are expected to have correct answers, these values should be specified as hidden fields in the items file; see discussion following (14).

¹⁵ If a line break is required within a single text field, use the HTML code

 .

¹⁶ If some items have more than one hidden field, it is important to make sure that each field is exclusively used for one purpose (e.g. field 5 for the correct answer to the item and field 6 for the expected answer to a comprehension question). It is possible to leave a field empty if a given item does not have a value associated with some field (e.g., no expected correct answer to an item, but there is an expected correct answer to the comprehension question leave field 5 empty, as a blank line, and input the expected correct answer to the comprehension question in field 6).

fields, the first specifying the sentence and the second, a picture. This corresponds to two substitution tags in the binary-image skeleton: {{field_1}} stands in for the sentence, and {{field_2}} for the URL in an tag, specifying the picture to be shown.

```
# target 1 indef-no.there
A boy is in the room.
http://DOMAIN/experimentpictures/definiteness-effect/picture1.png
# target 1 def-no.there
The boy is in the room.
http://DOMAIN/experimentpictures/definiteness-effect/picture1.png
# target 1 indef-there
There is a boy in the room
http://DOMAIN/experimentpictures/definiteness-effect/picture1.png
# target 1 def-there
There in the boy in the room
http://DOMAIN/experimentpictures/definiteness-effect/picture1.png
```

Figure 3. Sample item set with four conditions in *Lister*-compatible format

6. Lister: Creating randomized lists of items to be uploaded on Mechanical Turk

The Lister takes an items file as input and returns a comma-separated-value (.csv) format file with randomized lists of items that can be uploaded to AMT. When you run *lister.py* the program will maximally ask for the following information:¹⁸

- (15) a. the name of your items file (described in the previous section)
 - b. which sections should be treated as fillers
 - c. how many filler items you would like placed between each target item
 - d. how many filler items you would like placed in the beginning and end of the experiment
 - e. how many lists you would like to create
 - f. whether or not you would like the reverse of each list to also be created, to help reduce any ordering effects that may occur

¹⁷ Resources such as images and audio or video files must be hosted separately, rather than included with the AMT experiment. We recommend hosting such resources on your own web space, for example in a hidden directory on your website, and including links to those resources in the items.

18 We recommend that relevant scripts be copied into the same folder as your items file.

After this information is given, the *Lister* will create as many lists as requested, using a Latin Square design. That is, for each item set in a given section, an item for only one of the specified conditions for that item set will be included in a list. If the number of conditions does not match for all item sets in a section, an error will be given. Note that if all the item sets in a particular section have only one item, they will all be included, regardless of the condition names given.¹⁹

It is possible to dictate at least how many filler items must occur between each two target items and at the beginning and end of each list, if there are enough filler items. The algorithm creates counterbalanced and randomized lists for each section individually (including the fillers), and then spreads fillers between target items to minimally comply with the constraints given to the *Lister*. If additional filler items are left after this process, they are randomly distributed across the list. If there are not enough filler items to comply with the constraints given to the *Lister*, an error message will appear. Note that it is possible to have less filler items than target items, but in that case it is not possible to impose restrictions on their position relative to target items. It is also possible to have no sections designated as fillers. In this case, all the items in the file will be randomized first within their sections and then across other sections, but no other conditions can be imposed on the ordering.

If some items have more fields than others (for example, if some items have expected correct answers but others don't), the *Lister* will produce a message notifying the user of the discrepancy. The user may choose to go back and check that the file is correct or to continue with the current file. Assuming no such issues, the *Lister* will output a randomized lists file in CSV format (xxxxxxx.turk.csv) that can be uploaded onto AMT. For example, if the input to the *Lister* is the items file definiteness-effect.txt, the output will be definiteness-effect.turk.csv. The randomized lists file contains a row for each randomized list, indicating the actual text that the workers will see, with a header line at the top. AMT will substitute these strings into the HTML template for the participants to see.

7. Simulator: Verifying that the survey works properly

With your items randomized into lists and your HTML template prepared, you are in principle ready to upload your survey onto AMT. However, before doing so, we recommend pre-testing your study using the *simulator.py* script provided here. The *Simulator* takes as input a randomized lists file (xxxxxxx.turk.csv) and a template and creates a version of the survey based on one of the lists in the items file (you can choose which list to simulate). The *Simulator*

¹⁹ In a filler section where each item set is made up of one item each, it may be helpful to use different condition names even though all items will be displayed. This may help make the analysis script easier to create. For example, some fillers may be designated as 'catch' items and accuracy can be calculated for those items only and then used to discard participants from the analysis.

allows researchers to bypass the rather cumbersome process of uploading a survey onto AMT, or using AMT's *Sandbox* to simulate an experiment.²⁰ The output of the *Simulator* is an HTML file that can be opened in any web browser.

It is beneficial for researchers to complete their own study, for several reasons. Doing so helps spot any typos or errors in the rendering of the survey. If your items file includes hidden fields, you should make sure that all hidden fields are not displayed and all other fields are displayed as expected. Additionally, it is important to get a feeling for how hard the experiment is and how long it takes to complete it. This will help you determine payment for your HITs. Furthermore, because of the existence of *satisficing* behavior on AMT – the behavior of some workers to put in minimal effort and still get paid – it is useful to make sure that there are no easy strategies to complete the study without doing the linguistic task of interest. Strategies for dealing with such behavior include the addition of comprehension questions to items, the inclusion of 'catch' filler items that should have clear correct answers, and the use of diverse types of fillers that cannot all be solved in the same way.

8. Uploading surveys onto Mechanical Turk and downloading results files

After verifying that the survey is ready, it can be uploaded onto AMT. We refer the reader to our online instructions on uploading surveys onto AMT: https://github.com/mitcho/turktools/wiki. In some cases it may be necessary to reject submissions from workers who have completed more than one survey (see footnote 8 above). It may also be desirable to reject workers with low accuracy on comprehension questions and 'catch' filler items and those workers who failed to complete at least 80% of the study. To quickly identify workers who completed more than one survey, open the file in software such as Excel and sort by 'WorkerId.' Some suggestions and simple code for additional exclusion criteria and how to identify workers who meet them using R are provided in the analysis script analysis.r which is provided with *turktools* and can be modified by researchers to meet their own needs. In order to maintain a positive reputation as a requester on AMT, it is advisable to approve submissions in a timely fashion and to maintain a low rejection rate.

9. Decoder: Creating an output file in analyzable format

After your survey is completed, download a final copy of the raw results file from the AMT *Manage* tab. We recommend re-naming this file, e.g. definiteness-effect.results.csv

_

²⁰ Amazon Mechanical Turk's Sandbox allows researchers to upload "dummy" versions of their experiment, much like using the *Simulator*. However, to use the Sandbox it is necessary to set up the experiment independently on the Sandbox. This entire process must be repeated every time a mistake is discovered. Then one must go through all the same steps once again in order to create the actual survey on AMT; there is no way to transfer a survey from one platform to the other.

(to continue using the same example as above), and saving it in the same folder as the other documents produced in the course of preparing your survey. Use the *decoder.py* script, which should be in the same folder. The *Decoder* will ask for the name of the raw results file (here: definiteness-effect.results.csv) and return a decoded file, here named definiteness-effect.results.decoded.csv. This file can be analyzed in any commonly used statistical software.

Unlike the raw results file, the decoded results file contains one row for each *item* in each list in the survey. For each item, the decoded file specifies information about the Section, Condition, Item, List, PresentationOrder, and all fields and participant responses associated with that item, including hidden fields. In addition, the decoder adds a number of Factor columns, corresponding to different factor values. These factor values are constructed by simply splitting the Condition value up by hyphens. For example, in our example where an item may have Condition value 'indef-no.there', Factor1 would be set to 'indef' and Factor2 would be set to 'no.there.' This processing of condition names using the schema in (13) facilitates the analysis of the data later.

Metadata about a submission, including the AssignmentId (unique for each worker-survey pair), SubmissionStatus, WorkerId, WorkTimeInSeconds (for the entire survey) and answers to demographic questions, is duplicated for each item for a given worker in the file. Therefore, for each participant who worked on your survey there will be as many rows in the file as there were items in the survey,²² and the meta-information about the submission will be the same in all rows from that submission.

Once this process is complete, the results are ready to be analyzed by any standard statistics software. We have provided some basic code in the analysis script analysis.r, but researchers should adapt the script for their own needs.²³

10. Hosting AMT-style experiments on your own server with turkserver

In addition to posting an experiment on AMT, the tools we have provided with this paper can be used to generate surveys to be hosted on a researcher's own server. In order to do so, we have developed a program called *turkserver*, freely available at https://github.com/mitcho/turkserver

-

²¹ If condition names were not constructed using factor names and hyphens, as suggested in (13), custom code may be needed in the analysis script to recover factor values from the value of the Condition column.

²² But note that if a worker completed two surveys, for example, there will be twice as many rows. All the rows corresponding to the first survey completed by this worker will share meta-data pertaining to the first submission and all rows corresponding to the second survey will share meta-data pertaining to the second submission. The AssignmentIds will be distinct, but the WorkerId will be the same across these rows. This allows us to easily identify such workers who completed multiple surveys. At the time of writing, AMT cannot itself enforce that each worker only complete one survey per logical experiment.

²³ See also some analysis suggestions on the *turktools* website at https://github.com/mitcho/turktools/wiki.

under a liberal open-source license. *Turkserver* requires a server capable of serving PHP scripts, and is designed for Apache web servers. See the documentation at https://github.com/mitcho/turkserver for detailed instructions on installation and usage.

To use this option, construct the experiment as described in Sections 3-7. An experiment is set up with *turkserver* by uploading the template and item lists file (.turk.csv). A URL for the experiment is generated, and this URL can then be shared with potential participants. *Turkserver* handles the random assignment of participants to different lists, displaying the survey using the chosen list of items, and recording the results in a format compatible with AMT's results files.

There are some differences between running an experiment on AMT and on *turkserver*. AMT comes with many potential participants registered on the system, and therefore acts as a subject recruitment tool as well. On the other hand, the potential participant pool is limited to those people who have chosen to sign up to be a worker on AMT. With *turkserver* it is up to the experimenter to recruit subjects and send them to the experiment URL. Similarly, unlike AMT, *turkserver* does not handle payment to participants. *Turkserver* does, however, attempt to produce a unique AssignmentId for each submission, as AMT does, and to produce a WorkerId for each participant that is stable across experiments.

Running AMT-style surveys on *turkserver* can be useful for situations where the AMT subject pool does not include many speakers of a particular language, but such speakers can be recruited separately. It is also ideal for when the experimenter wishes to recruit participants from a particular group, such as a vetted subject pool, a class, or among colleagues. Furthermore, the ability to generate experiments that run both with and without AMT offers great flexibility. For example, one could use the exact same experimental materials (template and item list files) and conduct the experiment with the general public on AMT, and also on recruited non-naïve (working linguist) participants via *turkserver*, and compare the results in behavior.

11. Conclusion

We hope that the tools we have developed will allow linguists to quickly and efficiently test diverse hypotheses stemming from their work using linguistic surveys on Amazon's Mechanical Turk. The ease and speed with which linguistic hypotheses can be tested on AMT makes this platform a useful tool for theoretical linguists, and we hope that more linguists will learn to appreciate this tool and develop the skills necessary to use it. We hope that the tools we have developed will make experimental methods more accessible to researchers with less background in the use of such techniques and will thus make experimental methods more accessible to a wider range of working linguists.

References

Baayen, R. H. 2008. Analyzing linguistic data: a practical introduction to statistics using R. Cambridge, UK: Cambridge University Press.

Baayen, R.H., Davidson, D.J., Bates, D.M. 2008. Mixed-effects modeling with crossed random effects for subjects and items. Journal of Memory and Language, 59, 390-412.

Bard, E. G., Robertson, D. and Sorace, A. 1996. Magnitude estimation of linguistic acceptability. Language 72: 107–50.

Bates, D.M. and Maechler, M. 2009. lme4: Linear mixed-effects models using S4 classes. R package version 0.999375-32.

Barr, D.J., Levy, R., Scheepers, C., and Tily, H.J., 2013. Random effects structure for confirmatory hypothesis testing: Keep it maximal. Journal of Memory and Language 68(3):255-278.

Cable, S. and Harris, J. 2011. On the Grammatical Status of PP-Pied-Piping in English: Results from Sentence-Rating Experiments. Ms, ling.auf.net/lingbuzz/001302/current.pdf.

Chomsky, N. 1965. Aspects of the theory of syntax. Cambridge: MIT Press.

Cowart, W. 1997. Experimental syntax: applying objective methods to sentence judgments. Thousand Oaks, CA: Sage Publications.

Culicover, P. W., and Jackendoff, R. 2010. Quantitative methods alone are not enough: Response to Gibson and Fedorenko. Trends in Cognitive Sciences, 14(6), 234-235

DeGroot, M. H., Schervish, M. J., Fang, X., Lu, L., and Li, D. 1986. Probability and statistics. Reading, MA: Addison-Wesley.

Edelman, S., and Christiansen, M. 2003. How seriously should we take Minimalist syntax? Trends in Cognitive Sciences 7. 60–1.

Featherston, S. 2005. Magnitude estimation and what it can do for your syntax: some wh-constraints in German. Lingua 115. 1525–50.

Ferreira, F. 2005. Psycholinguistics, formal grammars, and cognitive science. The Linguistic Review 22. 365–80.

Fort, K., Adda, G., and Cohen, K.B. 2011. Amazon Mechanical Turk: Gold Mine or Coal Mine?. Computational Linguistics 37, No. 2, Pages 413-20.

Gelman, A., and Hill, J. 2007. Data analysis using regression and multilevel / hierarchical models. Cambridge, UK: Cambridge University Press.

Gibson, E., and Fedorenko, E. 2010. Weak quantitative standards in linguistics research. Trends in Cognitive Science 14. 233–4.

Gibson, E., Piantadosi, S., and Fedorenko, K. 2011. Using Mechanical Turk to Obtain and Analyze English Acceptability Judgments. Language and Linguistics Compass 5/8: 509–524.

Huang, Y.T., Spelke, E. and Snedeker, J. ms. What exactly do numbers mean?

Keller, F. 2000. Gradience in Grammar: Experimental and Computational Aspects of Degrees of Grammaticality. PhD Thesis, University of Edinburgh.

Marantz, A. 2005. Generative linguistics within the cognitive neuroscience of language. The Linguistic Review 22. 429–45.

Milsark, G. 1974. Existential sentences in English. Doctoral dissertation, MIT.

Milsark, G. 1977. Toward an explanation of certain peculiarities of the existential construction in English. Linguistic Analysis 3, 1-29.

Munro, R., Bethard, S., Kuperman V., Tzuyin Lai, V., Melnick, R., Potts, C., Schnoebelen, T., and Tily, H. 2010. Crowdsourcing and language studies: the new generation of linguistic data. Workshop on Creating Speech and Language Data with Amazon's Mechanical Turk, Los Angeles, CA.

Myers, J. 2009a. Syntactic judgment experiments. Language and Linguistics Compass 3. 406–23.

Myers, J. 2009b. The design and analysis of small-scale syntactic judgment experiments. Lingua 119. 425–44.

Pearson, H., Khan, M. and Snedeker, J. 2010. Even more evidence for the emptiness of plurality: An experimental investigation of plural interpretation as a species of implicature. Proceedings of SALT, 20, 489-508.

Phillips, C., and Lasnik, H. 2003. Linguistics and empirical evidence: Reply to Edelman and Christiansen. Trends in Cognitive Sciences, 7, 61-62.

Schütze, C. 1996. The empirical base of linguistics: Grammaticality judgments and linguistic methodology. Chicago: University of Chicago Press.

Schütze, C. and Sprouse, J. To appear. *Judgment Data. Research Methods in Linguistics*. Edited by Devyani Sharma and Rob Podesva.

Sprouse, J. 2009. Revisiting satiation: Evidence for an equalization response strategy. Linguistic Inquiry, 40, 329–341.

Sprouse, J. 2011. A validation of Amazon Mechanical Turk for the collection of acceptability judgments in linguistic theory. Behavior Research Methods 43.

Sprouse, J and Almeida, D. 2010. A quantitative defense of linguistic methodology. Irvine: University of California, ms.

Sprouse, J and Almeida, D. To appear. The empirical status of data in syntax: A reply to Gibson and Fedorenko. Language and Cognitive Processes.

Sprouse, J and Almeida, D. 2012. Assessing the reliability of textbook data in syntax: Adger's Core Syntax. Journal of Linguistics 48: 609-652.

Sprouse, J., Schütze, C., and Almeida, D. Ms. A comparison of informal and formal acceptability judgments using a random sample from Linguistic Inquiry 2001-2010.

Wasow, T., and Arnold, J. 2005. Intuitions in linguistic argumentation. Lingua 115. 1481–96.

Weskott, T., & Fanselow, G. 2011. On the Informativity of Different Measures of Linguistic Acceptability. Language, 87, 249-273.