# Mapping Natural Language Instructions to Mobile UI Action Sequences

**Yang Li    Jiacong He    Xin Zhou    Yuan Zhang    Jason Baldridge**

Google Research, Mountain View, CA, 94043

{liyang,zhouxin,zhangyua,jasonbaldridge}@google.com

## Abstract

We present a new problem: grounding natural language instructions to mobile user interface actions, and create three new datasets for it. For full task evaluation, we create PIXELHELP, a corpus that pairs English instructions with actions performed by people on a mobile UI emulator. To scale training, we decouple the language and action data by (a) annotating action phrase spans in HowTo instructions and (b) synthesizing grounded descriptions of actions for mobile user interfaces. We use a Transformer to extract action phrase tuples from long-range natural language instructions. A grounding Transformer then contextually represents UI objects using both their content and screen position and connects them to object descriptions. Given a starting screen and instruction, our model achieves 70.59% accuracy on predicting *complete* ground-truth action sequences in PIXELHELP.

## 1 Introduction

Language helps us work together to get things done. People instruct one another to coordinate joint efforts and accomplish tasks involving complex sequences of actions. This takes advantage of the abilities of different members of a speech community, e.g. a child asking a parent for a cup she cannot reach, or a visually impaired individual asking for assistance from a friend. Building computational agents able to help in such interactions is an important goal that requires true language grounding in environments where action matters.

An important area of language grounding involves tasks like completion of multi-step actions in a graphical user interface conditioned on language instructions (Branavan et al., 2009, 2010; Liu et al., 2018; Gur et al., 2019). These domains matter for accessibility, where language interfaces could help visually impaired individuals perform tasks with
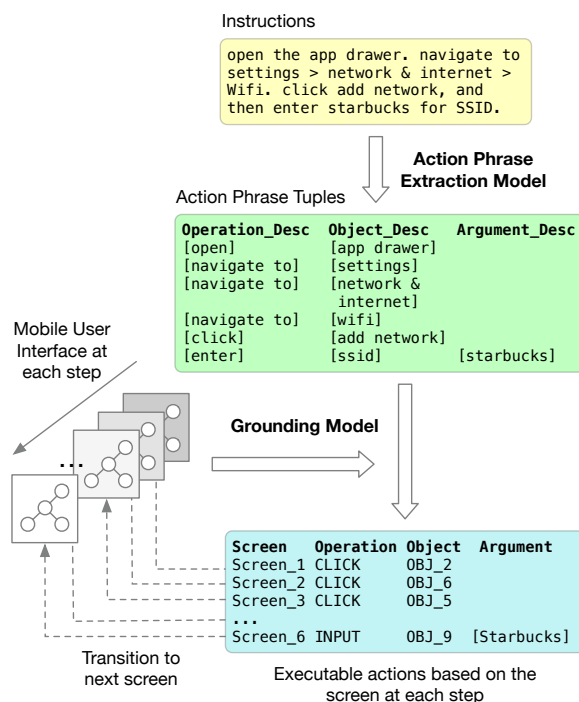


Figure 1: Our model extracts the phrase tuple that describe each action, including its operation, object and additional arguments, and grounds these tuples as executable action sequences in the UI.

interfaces that are predicated on sight. This also matters for *situational impairment* (Sarsenbayeva, 2018) when one cannot access a device easily while encumbered by other factors, such as cooking.

We focus on a new domain of task automation in which natural language instructions must be interpreted as a sequence of actions on a mobile touchscreen UI. Existing web search is quite capable of retrieving multi-step natural language instructions for user queries, such as "How to turn on flight mode on Android." Crucially, the missing piece for fulfilling the task automatically is to map the returned instruction to a sequence of actions that can be automatically executed on the device with

little user intervention; this our goal in this paper. This task automation scenario does not require a user to maneuver through UI details, which is useful for average users and is especially valuable for visually or situationally impaired users. The ability to execute an instruction can also be useful for other scenarios such as automatically examining the quality of an instruction.

Our approach (Figure 1) decomposes the problem into an *action phrase-extraction* step and a *grounding* step. The former extracts operation, object and argument descriptions from multi-step instructions; for this, we use Transformers (Vaswani et al., 2017) and test three span representations. The latter matches extracted operation and object descriptions with a UI object on a screen; for this, we use a Transformer that contextually represents UI objects and grounds object descriptions to them.

We construct three new datasets [1]. To assess full task performance on *naturally occurring* instructions, we create a dataset of 187 multi-step English instructions for operating Pixel Phones and produce their corresponding action-screen sequences using annotators. For action phrase extraction training and evaluation, we obtain English How-To instructions from the web and annotate action description spans. A Transformer with spans represented by sum pooling (Li et al., 2019) obtains 85.56% accuracy for predicting span sequences that completely match the ground truth. To train the grounding model, we synthetically generate 295k single-step commands to UI actions, covering 178K different UI objects across 25K mobile UI screens.

Our phrase extractor and grounding model together obtain 89.21% partial and 70.59% complete accuracy for matching ground-truth action sequences on this challenging task. We also evaluate alternative methods and representations of objects and spans and present qualitative analyses to provide insights into the problem and models.

## 2 Problem Formulation

Given an instruction of a multi-step task, $I = t_{1:n} = (t_1, t_2, ..., t_n)$, where $t_i$ is the $i$th token in instruction $I$, we want to generate a sequence of automatically executable actions, $a_{1:m}$, over a sequence of user interface screens $S$, with initial screen $s_1$

and screen transition function $s_j = \tau(a_{j-1}, s_{j-1})$:

$$p(a_{1:m}|s_1, \tau, t_{1:n}) = \prod_{j=1}^{m} p(a_j|a_{<j}, s_1, \tau, t_{1:n})$$
(1)

An action $a_j = [r_j, o_j, u_j]$ consists of an operation $r_j$ (e.g. `Tap` or `Text`), the UI object $o_j$ that $r_j$ is performed on (e.g., a button or an icon), and an additional argument $u_j$ needed for $o_j$ (e.g. the message entered in the chat box for `Text` or `null` for operations such as `Tap`). Starting from $s_1$, executing a sequence of actions $a_{<j}$ arrives at screen $s_j$ that represents the screen at the $j$th step: $s_j = \tau(a_{j-1}, \tau(...\tau(a_1, s_1)))$:

$$p(a_{1:m}|s_1, \tau, t_{1:n}) = \prod_{j=1}^{m} p(a_j|s_j, t_{1:n}) \quad (2)$$

Each screen $s_j = [c_{j,1:|s_j|}, \lambda_j]$ contains a set of UI objects and their structural relationships. $c_{j,1:|s_j|} = \{c_{j,k} \mid 1 \leq k \leq |s_j|\}$, where $|s_j|$ is the number of objects in $s_j$, from which $o_j$ is chosen. $\lambda_j$ defines the structural relationship between the objects. This is often a tree structure such as the *View* hierarchy for an Android interface[2] (similar to a DOM tree for web pages).

An instruction $I$ describes (possibly multiple) actions. Let $\bar{a}_j$ denote the phrases in $I$ that describes action $a_j$. $\bar{a}_j = [\bar{r}_j, \bar{o}_j, \bar{u}_j]$ represents a tuple of descriptions with each corresponding to a span—a subsequence of tokens—in $I$. Accordingly, $\bar{a}_{1:m}$ represents the description tuple sequence that we refer to as $\bar{a}$ for brevity. We also define $\bar{A}$ as all possible description tuple sequences of $I$, thus $\bar{a} \in \bar{A}$.

$$p(a_j|s_j, t_{1:n}) = \sum_{\bar{A}} p(a_j|\bar{a}, s_j, t_{1:n})p(\bar{a}|s_j, t_{1:n})$$
(3)

Because $a_j$ is independent of the rest of the instruction given its current screen $s_j$ and description $\bar{a}_j$, and $\bar{a}$ is only related to the instruction $t_{1:n}$, we can simplify (3) as (4).

$$p(a_j|s_j, t_{1:n}) = \sum_{\bar{A}} p(a_j|\bar{a}_j, s_j)p(\bar{a}|t_{1:n}) \quad (4)$$

We define $\hat{a}$ as the most likely description of actions for $t_{1:n}$.

$$\hat{a} = \arg\max_{\bar{a}} p(\bar{a}|t_{1:n})$$
$$= \arg\max_{\bar{a}_{1:m}} \prod_{j=1}^{m} p(\bar{a}_j|\bar{a}_{<j}, t_{1:n}) \quad (5)$$

This defines the action phrase-extraction model, which is then used by the grounding model:

$$p(a_j|s_j, t_{1:n}) \approx p(a_j|\hat{a}_j, s_j)p(\hat{a}_j|\hat{a}_{<j}, t_{1:n}) \quad (6)$$

$$p(a_{1:m}|t_{1:n}, S) \approx \prod_{j=1}^{m} p(a_j|\hat{a}_j, s_j)p(\hat{a}_j|\hat{a}_{<j}, t_{1:n}) \quad (7)$$

$p(\hat{a}_j|\hat{a}_{<j}, t_{1:n})$ identifies the description tuples for each action. $p(a_j|\hat{a}_j, s_j)$ grounds each description to an executable action given the screen.

## 3 Data

The ideal dataset would have natural instructions that have been executed by people using the UI. Such data can be collected by having annotators perform tasks according to instructions on a mobile platform, but this is difficult to scale. It requires significant investment to instrument: different versions of apps have different presentation and behaviors, and apps must be installed and configured for each task. Due to this, we create a small dataset of this form, PIXELHELP, for full task evaluation. For model training at scale, we create two other datasets: ANDROIDHOWTO for action phrase extraction and RICOSCA for grounding. Our datasets are targeted for English. We hope that starting with a high-resource language will pave the way to creating similar capabilities for other languages.

### 3.1 PIXELHELP Dataset

Pixel Phone Help pages[3] provide instructions for performing common tasks on Google Pixel phones such as *switch Wi-Fi settings* (Fig. 2) or *check emails*. Help pages can contain multiple tasks, with each task consisting of a sequence of steps. We pulled instructions from the help pages and kept ones that can be automatically executed. Instructions that requires additional user input such as *Tap the app you want to uninstall* are discarded.
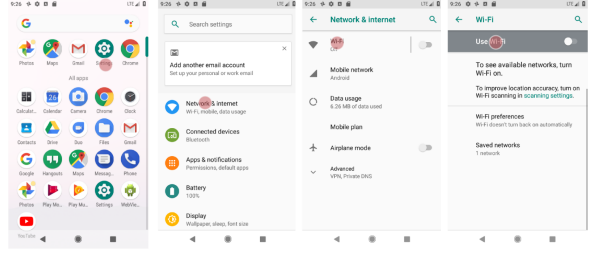
---

[3] https://support.google.com/pixelphone



Figure 2: PIXELHELP example: *Open your device's Settings app. Tap Network & internet. Click Wi-Fi. Turn on Wi-Fi.*. The instruction is paired with actions, each of which is shown as a red dot on a specific screen.

Also, instructions that involve actions on a physical button such as *Press the Power button for a few seconds* are excluded because these events cannot be executed on mobile platform emulators.

We instrumented a logging mechanism on a Pixel Phone emulator and had human annotators perform each task on the emulator by following the full instruction. The logger records every user action, including the type of touch events that are triggered, each object being manipulated, and screen information such as view hierarchies. Each item thus includes the instruction input, $t_{1:n}$, the screen for each step of task, $s_{1:m}$, and the target action performed on each screen, $a_{1:m}$.

In total, PIXELHELP includes 187 multi-step instructions of 4 task categories: 88 general tasks, such as configuring accounts, 38 Gmail tasks, 31 Chrome tasks, and 30 Photos related tasks. The number of steps ranges from two to eight, with a median of four. Because it has both natural instructions and grounded actions, we reserve PIXELHELP for evaluating full task performance.

### 3.2 ANDROIDHOWTO Dataset

No datasets exist that support learning the action phrase extraction model, $p(\hat{a}_j|\hat{a}_{<j}, t_{1:n})$, for mobile UIs. To address this, we extracted English instructions for operating Android devices by processing web pages to identify candidate instructions for how-to questions such as *how to change the input method for Android*. A web crawling service scrapes instruction-like content from various websites. We then filter the web contents using both heuristics and manual screening by annotators.

Annotators identified phrases in each instruction that describe executable actions. They were given a tutorial on the task and were instructed to skip instructions that are difficult to understand or label. For each component in an action description, they

select the span of words that describes the component using a web annotation interface (details are provided in the appendix). The interface records the start and end positions of each marked span. Each instruction was labeled by three annotators: three annotators agreed on 31% of full instructions and at least two agreed on 84%. For the consistency at the tuple level, the agreement across all the annotators is 83.6% for operation phrases, 72.07% for object phrases, and 83.43% for input phrases. The discrepancies are usually small, e.g., a description marked as *your Gmail address* or *Gmail address*.

The final dataset includes 32,436 data points from 9,893 unique How-To instructions and split into training (8K), validation (1K) and test (900). All test examples have perfect agreement across all three annotators for the *entire* sequence. In total, there are 190K operation spans, 172K object spans, and 321 input spans labeled. The lengths of the instructions range from 19 to 85 tokens, with median of 59. They describe a sequence of actions from one to 19 steps, with a median of 5.

### 3.3 RICOSCA Dataset

Training the grounding model, $p(a_j|\hat{a}_j, s_j)$ involves pairing action tuples $a_j$ along screens $s_j$ with action description $\hat{a}_j$. It is very difficult to collect such data at scale. To get past the bottleneck, we exploit two properties of the task to generate a *synthetic* command-action dataset, RICOSCA. First, we have precise structured and visual knowledge of the UI layout, so we can spatially relate UI elements to each other and the overall screen. Second, a grammar grounded in the UI can cover many of the commands and kinds of reference needed for the problem. This does not capture all manners of interacting conversationally with a UI, but it proves effective for training the grounding model.

Rico is a public UI corpus with 72K Android UI screens mined from 9.7K Android apps (Deka et al., 2017). Each screen in Rico comes with a screenshot image and a view hierarchy of a collection of UI objects. Each individual object, $c_{j,k}$, has a set of properties, including its name (often an English phrase such as *Send*), type (e.g., `Button`, `Image` or `Checkbox`), and bounding box position on the screen. We manually removed screens whose view hierarchies do not match their screenshots by asking annotators to visually verify whether the bounding boxes of view hierarchy leaves match each UI object on the corresponding screenshot image. This

filtering results in 25K unique screens.

For each screen, we randomly select UI elements as target objects and synthesize commands for operating them. We generate multiple commands to capture different expressions describing the operation $\hat{r}_j$ and the target object $\hat{o}_j$. For example, the `Tap` operation can be referred to as *tap*, *click*, or *press*. The template for referring to a target object has slots `Name`, `Type`, and `Location`, which are instantiated using the following strategies:

- *Name-Type*: the target's name and/or type (*the OK button* or *OK*).
- *Absolute-Location*: the target's screen location (*the menu at the top right corner*).
- *Relative-Location*: the target's relative location to other objects (*the icon to the right of Send*).

Because all commands are synthesized, the span that describes each part of an action, $\hat{a}_j$ with respect to $t_{1:n}$, is known. Meanwhile, $a_j$ and $s_j$, the actual action and the associated screen, are present because the constituents of the action are synthesized. In total, RICOSCA contains 295,476 single-step synthetic commands for operating 177,962 different target objects across 25,677 Android screens.

## 4 Model Architectures

Equation 7 has two parts. $p(\hat{a}_j|\hat{a}_{<j}, t_{1:n})$ finds the best phrase tuple that describes the action at the $j$th step given the instruction token sequence. $p(a_j|\hat{a}_j, s_j)$ computes the probability of an executable action $a_j$ given the best description of the action, $\hat{a}_j$, and the screen $s_j$ for the $j$th step.

### 4.1 Phrase Tuple Extraction Model

A common choice for modeling the conditional probability $p(\bar{a}_j|\bar{a}_{<j}, t_{1:n})$ (see Equation 5) are encoder-decoders such as LSTMs (Hochreiter and Schmidhuber, 1997) and Transformers (Vaswani et al., 2017). The output of our model corresponds to positions in the input sequence, so our architecture is closely related to Pointer Networks (Vinyals et al., 2015).

Figure 3 depicts our model. An encoder $g$ computes a latent representation $h_{1:n} \in R^{n \times |h|}$ of the tokens from their embeddings: $h_{1:n} = g(e(t_{1:n}))$. A decoder $f$ then generates the hidden state $q_j = f(q_{<j}, \bar{a}_{<j}, h_{1:n})$ which is used to compute a query vector that locates each phrase of a tuple $(\bar{r}_j, \bar{o}_j, \bar{u}_j)$ at each step. $\bar{a}_j = [\bar{r}_j, \bar{o}_j, \bar{u}_j]$ and they
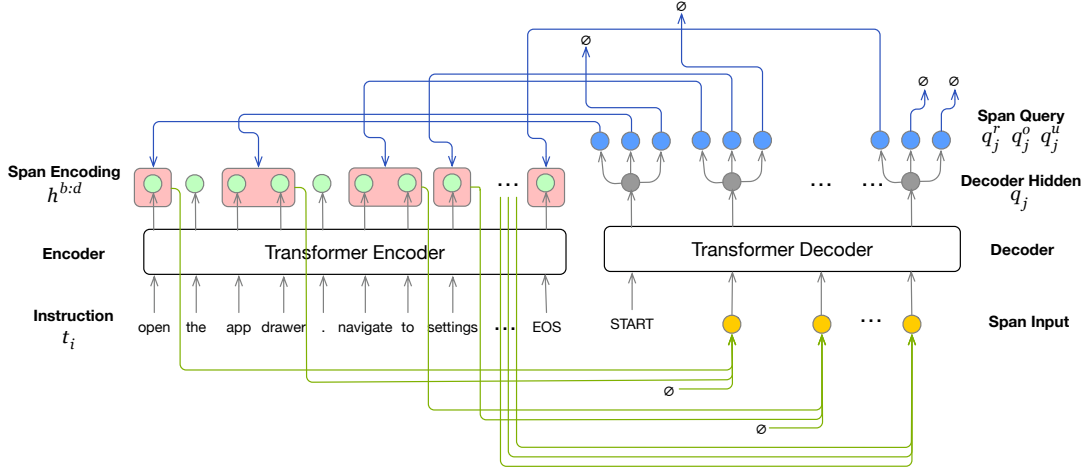
Figure 3: The Phrase Tuple Extraction model encodes the instruction's token sequence and then outputs a tuple sequence by querying into all possible spans of the encoded sequence. Each tuple contains the span positions of three phrases in the instruction that describe the action's operation, object and optional arguments, respectively, at each step. $\varnothing$ indicates the phrase is missing in the instruction and is represented by a special span encoding.

are assumed conditionally independent given previously extracted phrase tuples and the instruction, so $p(\bar{a}_j|\bar{a}_{<j}, t_{1:n}) = \prod_{\bar{y} \in \{\bar{r}, \bar{o}, \bar{u}\}} p(\bar{y}_j|\bar{a}_{<j}, t_{1:n})$.

Note that $\bar{y}_j \in \{\bar{r}_j, \bar{o}_j, \bar{u}_j\}$ denotes a specific span for $y \in \{r, o, u\}$ in the action tuple at step $j$. We therefore rewrite $\bar{y}_j$ as $y_j^{b:d}$ to explicitly indicate that it corresponds to the span for $r$, $o$ or $u$, starting at the $b$th position and ending at the $d$th position in the instruction, $1 \le b < d \le n$. We now parameterize the conditional probability as:

$$p(y_j^{b:d}|\bar{a}_{<j}, t_{1:n}) = \text{softmax}(\alpha(q_j^y, h^{b:d}))$$
$$y \in \{r, o, u\} \tag{8}$$

As shown in Figure 3, $q_j^y$ indicates task-specific query vectors for $y \in \{r, o, u\}$. They are computed as $q_j^y = \phi(q_j, \theta_y)W_y$, a multi-layer perceptron followed by a linear transformation. $\theta_y$ and $W_y$ are trainable parameters. We use separate parameters for each of $r$, $o$ and $u$. $W_y \in R^{|\phi_y| \times |h|}$ where $|\phi_y|$ is the output dimension of the multi-layer perceptron. The alignment function $\alpha(\cdot)$ scores how a query vector $q_j^y$ matches a span whose vector representation $h^{b:d}$ is computed from encodings $h_{b:d}$.

**Span Representation.** There are a quadratic number of possible spans given a token sequence (Lee et al., 2017), so it is important to design a fixed-length representation $h^{b:d}$ of a variable-length token span that can be *quickly* computed. Beginning-Inside-Outside (BIO) (Ramshaw and Marcus, 1995)–commonly used to indicate spans in tasks such as named entity recognition–marks

whether each token is beginning, inside, or outside a span. However, BIO is not ideal for our task because subsequences for describing different actions can overlap, e.g., in *click X and Y*, *click* participates in both actions *click X* and *click Y*. In our experiments we consider several recent, more flexible span representations (Lee et al., 2016, 2017; Li et al., 2019) and show their impact in Section 5.2.

With fixed-length span representations, we can use common alignment techniques in neural networks (Bahdanau et al., 2014; Luong et al., 2015). We use the dot product between the query vector and the span representation: $\alpha(q_j^y, h^{b:d}) = q_j^y \cdot h^{b:d}$ At each step of decoding, we feed the previously decoded phrase tuples, $\bar{a}_{<j}$ into the decoder. We can use the concatenation of the vector representations of the three elements in a phrase tuple or the sum their vector representations as the input for each decoding step. The entire phrase tuple extraction model is trained by minimizing the softmax cross entropy loss between the predicted and ground-truth spans of a sequence of phrase tuples.

## 4.2 Grounding Model

Having computed the sequence of tuples that best describe each action, we connect them to executable actions based on the screen at each step with our grounding model (Fig. 4). In step-by-step instructions, each part of an action is often clearly stated. Thus, we assume the probabilities of the operation $r_j$, object $o_j$, and argument $u_j$ are
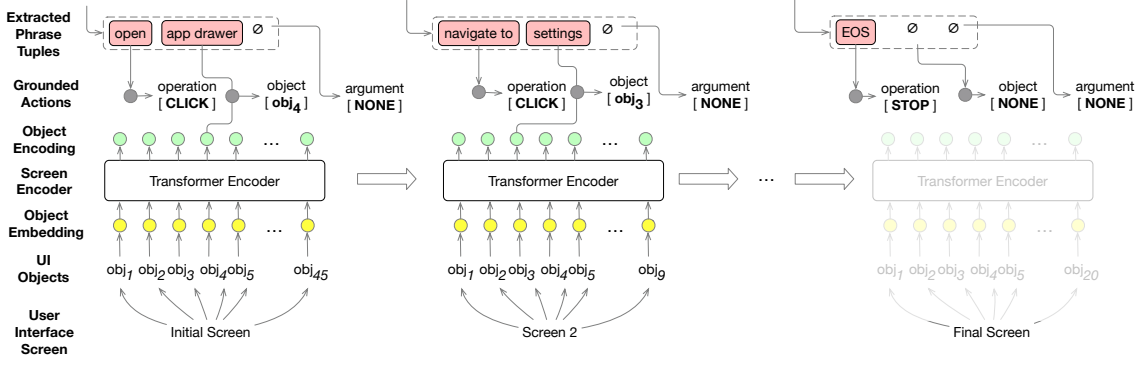
Figure 4: The Grounding model grounds each phrase tuple extracted by the Phrase Extraction model as an operation type, a screen-specific object ID, and an argument if present, based on a contextual representation of UI objects for the given screen. A grounded action tuple can be automatically executed.

independent given their description and the screen.

$$p(a_j|\hat{a}_j, s_j) = p([r_j, o_j, u_j]|[\hat{r}_j, \hat{o}_j, \hat{u}_j], s_j)$$
$$= p(r_j|\hat{r}_j, s_j)p(o_j|\hat{o}_j, s_j)p(u_j|\hat{u}_j, s_j)$$
$$= p(r_j|\hat{r}_j)p(o_j|\hat{o}_j, s_j)$$
$$(9)$$

We simplify with two assumptions: (1) an operation is often fully described by its instruction without relying on the screen information and (2) in mobile interaction tasks, an argument is only present for the Text operation, so $u_j = \hat{u}_j$. We parameterize $p(r_j|\hat{r}_j)$ as a feedforward neural network:

$$p(r_j|\hat{r}_j) = \text{softmax}(\phi(\hat{r}'_j, \theta_r)W_r) \qquad (10)$$

$\phi(\cdot)$ is a multi-layer perceptron with trainable parameters $\theta_r$. $W^r \in R^{|\phi_r| \times |r|}$ is also trainable, where $|\phi_r|$ is the output dimension of the $\phi(\cdot, \theta_r)$ and $|r|$ is the vocabulary size of the operations. $\phi(\cdot)$ takes the sum of the embedding vectors of each token in the operation description $\hat{r}_j$ as the input: $\hat{r}'_j = \sum_{k=b}^{d} e(t_k)$ where $b$ and $d$ are the start and end positions of $\hat{r}_j$ in the instruction.

Determining $o_j$ is to select a UI object from a variable-number of objects on the screen, $c_{j,k} \in s_j$ where $1 \le k \le |s_j|$, based on the given object description, $\hat{o}_j$. We parameterize the conditional probability as a deep neural network with a softmax output layer taking logits from an alignment function:

$$p(o_j|\hat{o}_j, s_j) = p(o_j = c_{j,k}|\hat{o}_j, c_{j,1:|s_j|}, \lambda_j)$$
$$= \text{softmax}(\alpha(\hat{o}'_j, c'_{j,k})) \qquad (11)$$

The alignment function $\alpha(\cdot)$ scores how the object description vector $\hat{o}'_j$ matches the latent representation of each UI object, $c'_{j,k}$. This can be as simple as the dot product of the two vectors. The latent representation $\hat{o}'_j$ is acquired with a multi-layer perceptron followed by a linear projection:

$$\hat{o}'_j = \phi(\sum_{k=b}^{d} e(t_k), \theta_o)W_o \qquad (12)$$

$b$ and $d$ are the start and end index of the object description $\hat{o}_j$. $\theta_o$ and $W_o$ are trainable parameters with $W_o \in R^{|\phi_o| \times |o|}$, where $|\phi_o|$ is the output dimension of $\phi(\cdot, \theta_o)$ and $|o|$ is the dimension of the latent representation of the object description.

**Contextual Representation of UI Objects.** To compute latent representations of each candidate object, $c'_{j,k}$, we use both the object's properties and its context, i.e., the structural relationship with other objects on the screen. There are different ways for encoding a variable-sized collection of items that are structurally related to each other, including Graph Convolutional Networks (GCN) (Niepert et al., 2016) and Transformers (Vaswani et al., 2017). GCNs use an adjacency matrix predetermined by the UI structure to regulate how the latent representation of an object should be affected by its neighbors. Transformers allow each object to carry its own positional encoding, and the relationship between objects can be learned instead.

The input to the Transformer encoder is a combination of the *content embedding* and the *positional encoding* of each object. The content properties of an object include its name and type. We compute the content embedding of by concatenating the name embedding, which is the average embedding of the bag of tokens in the object name, and the

type embedding. The positional properties of an object include both its spatial position and structural position. The spatial positions include the top, left, right and bottom screen coordinates of the object. We treat each of these coordinates as a discrete value and represent it via an embedding. Such a feature representation for coordinates was used in ImageTransformer to represent pixel positions in an image (Parmar et al., 2018). The spatial embedding of the object is the sum of these four coordinate embeddings. To encode structural information, we use the index positions of the object in the preorder and the postorder traversal of the view hierarchy tree, and represent these index positions as embeddings in a similar way as representing coordinates. The content embedding is then summed with positional encodings to form the embedding of each object. We then feed these object embeddings into a Transformer encoder model to compute the latent representation of each object, $c'_{j,k}$.

The grounding model is trained by minimizing the cross entropy loss between the predicted and ground-truth object and the loss between the predicted and ground-truth operation.

# 5   Experiments

Our goal is to develop models and datasets to map multi-step instructions into automatically executable actions given the screen information. As such, we use PIXELHELP's paired natural instructions and action-screen sequences solely for testing. In addition, we investigate the model quality on phrase tuple extraction tasks, which is a crucial building block for the overall grounding quality[4].

## 5.1   Datasets and Metrics

We use two metrics that measure how a predicted tuple sequence matches the ground-truth sequence.
- *Complete Match*: The score is $1$ if two sequences have the same length and have the identical tuple $[\hat{r}_j, \hat{o}_j, \hat{u}_j]$ at each step, otherwise $0$.
- *Partial Match*: The number of steps of the predicted sequence that match the ground-truth sequence divided by the length of the ground-truth sequence (ranging between $0$ and $1$).

We train and validate using ANDROIDHOWTO and RICOSCA, and evaluate on PIXELHELP. During training, single-step synthetic command-action

---

[4]Our model code is released at `https : / / github . com / google-research / google-research / tree/master/seq2act`.

| Span Rep. $h^{b:d}$ | Partial | Complete |
|---|---|---|
| Sum Pooling $\sum_{k=b}^{d} h_k$ | 92.80 | 85.56 |
| StartEnd Concat$[h_b; h_d]$ | 91.94 | 84.56 |
| $[h_b; h_d, \hat{e}^{b:d}, \phi(d-b)]$ | 91.11 | 84.33 |

Table 1: ANDROIDHOWTO phrase tuple extraction test results using different span representations $h^{b:d}$ in (8). $\hat{e}^{b:d} = \sum_{k=b}^{d} w(h_k)e(t_k)$, where $w(\cdot)$ is a learned weight function for each token embedding (Lee et al., 2017). See the pseudocode for fast computation of these in the appendix.

examples are dynamically stitched to form sequence examples with a certain length distribution. To evaluate the full task, we use Complete and Partial Match on grounded action sequences $a_{1:m}$ where $a_j = [r_j, o_j, u_j]$.

The token vocabulary size is 59K, which is compiled from both the instruction corpus and the UI name corpus. There are 15 UI types, including 14 common UI object types, and a type to catch all less common ones. The output vocabulary for operations include CLICK, TEXT, SWIPE and EOS.

## 5.2   Model Configurations and Results

**Tuple Extraction.** For the action-tuple extraction task, we use a 6-layer Transformer for both the encoder and the decoder. We evaluate three different span representations. Area Attention (Li et al., 2019) provides a parameter-free representation of each possible span (one-dimensional area), by summing up the encoding of each token in the subsequence: $h^{b:d} = \sum_{k=b}^{d} h_k$. The representation of each span can be computed in constant time invariant to the length of the span, using a summed area table. Previous work concatenated the encoding of the start and end tokens as the span representation, $h^{b:d} = [h_b; h_d]$ (Lee et al., 2016) and a generalized version of it (Lee et al., 2017). We evaluated these three options and implemented the representation in Lee et al. (2017) using a summed area table similar to the approach in area attention for fast computation. For hyperparameter tuning and training details, refer to the appendix.

Table 1 gives results on ANDROIDHOWTO's test set. All the span representations perform well. Encodings of each token from a Transformer already capture sufficient information about the entire sequence, so even only using the start and end encodings yields strong results. Nonetheless, area attention provides a small boost over the others. As a new dataset, there is also considerable headroom remaining, particularly for complete match.

| Screen Encoder | Partial | Complete |
|---|---|---|
| Heuristic | 62.44 | 42.25 |
| Filter-1 GCN | 76.44 | 52.41 |
| Distance GCN | 82.50 | 59.36 |
| Transformer | 89.21 | 70.59 |

Table 2: PIXELHELP grounding accuracy. The differences are statistically significant based on t-test over 5 runs ($p < 0.05$).

**Grounding.** For the grounding task, we compare Transformer-based screen encoder for generating object representations $h^{b:d}$ with two baseline methods based on graph convolutional networks. The *Heuristic* baseline matches extracted phrases against object names directly using BLEU scores. *Filter-1 GCN* performs graph convolution without using adjacent nodes (objects), so the representation of each object is computed only based on its own properties. *Distance GCN* uses the distance between objects in the view hierarchy, i.e., the number of edges to traverse from one object to another following the tree structure. This contrasts with the traditional GCN definition based on adjacency, but is needed because UI objects are often leaves in the tree; as such, they are not adjacent to each other structurally but instead are connected through non-terminal (container) nodes. Both Filter-1 GCN and Distance GCN use the same number of parameters (see the appendix for details).

To train the grounding model, we first train the Tuple Extraction sub-model on ANDROIDHOWTO and RICOSCA. For the latter, only language related features (commands and tuple positions in the command) are used in this stage, so screen and action features are not involved. We then freeze the Tuple Extraction sub-model and train the grounding sub-model on RICOSCA using both the command and screen-action related features. The screen token embeddings of the grounding sub-model share weights with the Tuple Extraction sub-model.

Table 2 gives full task performance on PIXEL-HELP. The Transformer screen encoder achieves the best result with 70.59% accuracy on Complete Match and 89.21% on Partial Match, which sets a strong baseline result for this new dataset while leaving considerable headroom. The GCN-based methods perform poorly, which shows the importance of contextual encodings of the information from other UI objects on the screen. Distance GCN does attempt to capture context for UI objects that

are structurally close; however, we suspect that the distance information that is derived from the view hierarchy tree is noisy because UI developers can construct the structure differently for the same UI.[5] As a result, the strong bias introduced by the structure distance does not always help. Nevertheless, these models still outperformed the heuristic baseline that achieved 62.44% for partial match and 42.25% for complete match.

### 5.3 Analysis

To explore how the model grounds an instruction on a screen, we analyze the relationship between words in the instruction language that refer to specific locations on the screen, and actual positions on the UI screen. We first extract the embedding weights from the trained phrase extraction model for words such as *top*, *bottom*, *left* and *right*. These words occur in object descriptions such as *the check box at the top of the screen*. We also extract the embedding weights of object screen positions, which are used to create object positional encoding. We then calculate the correlation between word embedding and screen position embedding using cosine similarity. Figure 5 visualizes the correlation as a heatmap, where brighter colors indicate higher correlation. The word *top* is strongly correlated with the top of the screen, but the trend for other location words is less clear. While *left* is strongly correlated with the left side of the screen, other regions on the screen also show high correlation. This is likely because *left* and *right* are not only used for referring to absolute locations on the screen, but also for relative spatial relationships, such as *the icon to the left of the button*. For *bottom*, the strongest correlation does not occur at the very bottom of the screen because many UI objects in our dataset do not fall in that region. The region is often reserved for system actions and the on-screen keyboard, which are not covered in our dataset.

The phrase extraction model passes phrase tuples to the grounding model. When phrase extraction is incorrect, it can be difficult for the grounding model to predict a correct action. One way to mitigate such cascading errors is using the hidden state of the phrase decoding model at each step, $q_j$. Intuitively, $q_j$ is computed with the access to the encoding of each token in the instruction via the Transformer encoder-decoder attention, which can

---

[5]While it is possible to directly use screen visual data for grounding, detecting UI objects from raw pixels is nontrivial. It would be ideal to use both structural and visual data.
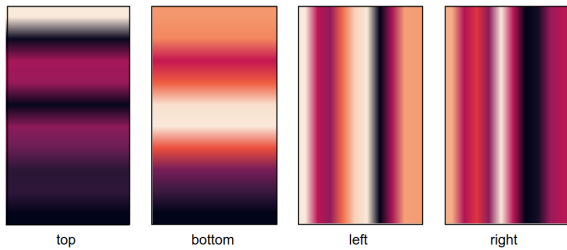
Figure 5: Correlation between location-related words in instructions and object screen position embedding.

potentially be a more robust span representation. However, in our early exploration, we found that grounding with $q_j$ performs stunningly well for grounding RICOSCA validation examples, but performs poorly on PIXELHELP. The learned hidden state likely captures characteristics in the synthetic instructions and action sequences that do not manifest in PIXELHELP. As such, using the hidden state to ground remains a challenge when learning from unpaired instruction-action data.

The phrase model failed to extract correct steps for 14 tasks in PIXELHELP. In particular, it resulted in extra steps for 11 tasks and extracted incorrect steps for 3 tasks, but did not skip steps for any tasks. These errors could be caused by different language styles manifested by the three datasets. Synthesized commands in RICOSCA tend to be brief. Instructions in ANDROIDHOWTO seem to give more contextual description and involve diverse language styles, while PIXELHELP often has a more consistent language style and gives concise description for each step.

## 6 Related Work

Previous work (Branavan et al., 2009, 2010; Liu et al., 2018; Gur et al., 2019) investigated approaches for grounding natural language on desktop or web interfaces. Manuvinakurike et al. (2018) contributed a dataset for mapping natural language instructions to actionable image editing commands in Adobe Photoshop. Our work focuses on a new domain of grounding natural language instructions into executable actions on mobile user interfaces. This requires addressing modeling challenges due to the lack of paired natural language and action data, which we supply by harvesting rich instruction data from the web and synthesizing UI commands based on a large scale Android corpus.

Our work is related to semantic parsing, particularly efforts for generating executable outputs such as SQL queries (Suhr et al., 2018). It is also broadly related to language grounding in the human-robot interaction literature where human dialog results in robot actions (Khayrallah et al., 2015).

Our task setting is closely related to work on language-conditioned navigation, where an agent executes an instruction as a sequence of movements (Chen and Mooney, 2011; Mei et al., 2016; Misra et al., 2017; Anderson et al., 2018; Chen et al., 2019). Operating user interfaces is similar to navigating the physical world in many ways. A mobile platform consists of millions of apps that each is implemented by different developers independently. Though platforms such as Android strive to achieve interoperability (e.g., using Intent or AIDL mechanisms), apps are more often than not built by convention and do not expose programmatic ways for communication. As such, each app is opaque to the outside world and the only way to manipulate it is through its GUIs. These hurdles while working with a vast array of existing apps are like physical obstacles that cannot be ignored and must be negotiated contextually in their given environment.

## 7 Conclusion

Our work provides an important first step on the challenging problem of grounding natural language instructions to mobile UI actions. Our decomposition of the problem means that progress on either can improve full task performance. For example, action span extraction is related to both semantic role labeling (He et al., 2018) and extraction of multiple facts from text (Jiang et al., 2019) and could benefit from innovations in span identification and multitask learning. Reinforcement learning that has been applied in previous grounding work may help improve out-of-sample prediction for grounding in UIs and improve direct grounding from hidden state representations. Lastly, our work provides a technical foundation for investigating user experiences in language-based human computer interaction.

# References

Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton van den Hengel. 2018. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473.

S. R. K. Branavan, Harr Chen, Luke S. Zettlemoyer, and Regina Barzilay. 2009. Reinforcement learning for mapping instructions to actions. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 1 - Volume 1*, ACL '09, pages 82–90, Stroudsburg, PA, USA. Association for Computational Linguistics.

S.R.K. Branavan, Luke Zettlemoyer, and Regina Barzilay. 2010. Reading between the lines: Learning to map high-level instructions to commands. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 1268–1277, Uppsala, Sweden. Association for Computational Linguistics.

David L. Chen and Raymond J. Mooney. 2011. Learning to interpret natural language navigation instructions from observations. In *Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence*, AAAI'11, pages 859–865. AAAI Press.

Howard Chen, Alane Suhr, Dipendra Misra, and Yoav Artzi. 2019. Touchdown: Natural language navigation and spatial reasoning in visual street environments. In *Conference on Computer Vision and Pattern Recognition*.

Biplab Deka, Zifeng Huang, Chad Franzen, Joshua Hibschman, Daniel Afergan, Yang Li, Jeffrey Nichols, and Ranjitha Kumar. 2017. Rico: A mobile app dataset for building data-driven design applications. In *Proceedings of the 30th Annual Symposium on User Interface Software and Technology*, UIST '17.

Izzeddin Gur, Ulrich Rueckert, Aleksandra Faust, and Dilek Hakkani-Tur. 2019. Learning to navigate the web. In *International Conference on Learning Representations*.

Luheng He, Kenton Lee, Omer Levy, and Luke Zettlemoyer. 2018. Jointly predicting predicates and arguments in neural semantic role labeling. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 364–369, Melbourne, Australia. Association for Computational Linguistics.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Comput.*, 9(8):1735–1780.

Tianwen Jiang, Tong Zhao, Bing Qin, Ting Liu, Nitesh Chawla, and Meng Jiang. 2019. Multi-input multi-output sequence labeling for joint extraction of fact and condition tuples from scientific text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 302–312, Hong Kong, China. Association for Computational Linguistics.

Huda Khayrallah, Sean Trott, and Jerome Feldman. 2015. Natural language for human robot interaction.

Kenton Lee, Luheng He, Mike Lewis, and Luke Zettlemoyer. 2017. End-to-end neural coreference resolution. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 188–197, Copenhagen, Denmark. Association for Computational Linguistics.

Kenton Lee, Tom Kwiatkowski, Ankur P. Parikh, and Dipanjan Das. 2016. Learning recurrent span representations for extractive question answering. *CoRR*, abs/1611.01436.

Yang Li, Lukasz Kaiser, Samy Bengio, and Si Si. 2019. Area attention. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 3846–3855, Long Beach, California, USA. PMLR.

E. Z. Liu, K. Guu, P. Pasupat, T. Shi, and P. Liang. 2018. Reinforcement learning on web interfaces using workflow-guided exploration. In *International Conference on Learning Representations (ICLR)*.

Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Effective approaches to attention-based neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1412–1421, Lisbon, Portugal. Association for Computational Linguistics.

Ramesh Manuvinakurike, Jacqueline Brixey, Trung Bui, Walter Chang, Doo Soon Kim, Ron Artstein, and Kallirroi Georgila. 2018. Edit me: A corpus and a framework for understanding natural language image editing. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC-2018)*, Miyazaki, Japan. European Languages Resources Association (ELRA).

Hongyuan Mei, Mohit Bansal, and Matthew R. Walter. 2016. Listen, attend, and walk: Neural mapping of navigational instructions to action sequences. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, AAAI'16, pages 2772–2778. AAAI Press.

Dipendra Misra, John Langford, and Yoav Artzi. 2017. Mapping instructions and visual observations to actions with reinforcement learning. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1004–1015.

Mathias Niepert, Mohamed Ahmed, and Konstantin Kutzkov. 2016. Learning convolutional neural networks for graphs. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2014–2023, New York, New York, USA. PMLR.

Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. 2018. Image transformer. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 4055–4064, Stockholmsmssan, Stockholm Sweden. PMLR.

Lance Ramshaw and Mitch Marcus. 1995. Text chunking using transformation-based learning. In *Third Workshop on Very Large Corpora*.

Zhanna Sarsenbayeva. 2018. Situational impairments during mobile interaction. In *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, pages 498–503.

Alane Suhr, Srinivasan Iyer, and Yoav Artzi. 2018. Learning to map context-dependent sentences to executable formal queries. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2238–2249, New Orleans, Louisiana. Association for Computational Linguistics.

Richard Szeliski. 2010. *Computer Vision: Algorithms and Applications*, 1st edition. Springer-Verlag, Berlin, Heidelberg.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *CoRR*, abs/1706.03762.

Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. 2015. Pointer networks. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 2692–2700. Curran Associates, Inc.

## A  Data

We present the additional details and analysis of the datasets. To label action phrase spans for the ANDROIDHOWTO dataset, 21 annotators (9 males and 12 females, 23 to 28 years old) were employed as contractors. They were paid hourly wages that are competitive for their locale. They have standard rights as contractors. They were native English speakers, and rated themselves 4 out of 5 regarding their familiarity with Android (1: not familiar and 5: very familiar).

Each annotator is presented a web interface, depicted in Figure 6. The instruction to be labeled is shown on the left of the interface. From the instruction, the annotator is asked to extract a sequence of action phrase tuples on the right, providing one tuple per row. Before a labeling session, an annotator is asked to go through the annotation guidelines, which are also accessible throughout the session.

To label each tuple, the annotator first indicates the type of operation (Action Type) the step is about by selecting from `Click`, `Swipe`, `Input` and `Others` (the catch-all category). The annotator then uses a mouse to select the phrase in the instruction for "Action Verb" (i.e., operation description) and for "object description". A selected phrase span is automatically shown in the corresponding box and the span positions in the instruction are recorded. If the step involves an additional argument, the annotator clicks on "Content Input" and then marks a phrase span in the instruction (see the second row). Once finished with creating a tuple, the annotator moves onto the next tuple by clicking the "+" button on the far right of the interface along the row, which inserts an empty tuple after the row. The annotator can delete a tuple (row) by clicking the "-" button on the row. Finally, the annotator clicks on the "Submit" button at the bottom of the screen to finish a session.

The lengths of the instructions range from 19 to 85 tokens, with median of 59, and they describe a sequence of actions from 1 to 19 steps, with a median of 5. Although the description for operations tend to be short (most of them are one to two words), the description for objects can vary dramatically in length, ranging from 1 to 19. The large range of description span lengths requires an efficient algorithm to compute its representation.

## B  Computing Span Representations

We evaluated three types of span representations. Here we give details on how each representation is computed. For sum pooling, we use the implementation of area attention (Li et al., 2019) that allows constant time computation of the representation of each span by using summed area tables. The TensorFlow implementation of the representation

Figure 6: The web interface for annotators to label action phrase spans in an ANDROIDHOWTO instruction.

is available on Github[6].

---

**Algorithm 1:** Compute the Start-End Concat span representation for all spans in parallel.

---

**Input:** A tensor $H$ in shape of $[L, D]$ that represents a sequence of vector with length $L$ and depth $D$.

**Output:** representation of each span, $U$.

1 **Hyperparameter**: max span width $M$.
2 **Init** start & end tensor: $S \leftarrow H, E \leftarrow H$;
3 **for** $m = 1, \cdots, M - 1$ **do**
4     $S' \leftarrow H[: -m, :]$ ;
5     $E' \leftarrow H[m :, :]$ ;
6     $S \leftarrow [S \ S']$, concat on the 1st dim;
7     $E \leftarrow [E \ E']$, concat on the 1st dim;
8 $U \leftarrow [S \ E]$, concat on the last dim;
9 **return** $U$.

---

Algorithm 1 gives the recipe for Start-End Concat (Lee et al., 2016) using Tensor operations. The advanced form (Lee et al., 2017) takes two other features: the weighted sum over all the token embedding vectors within each span and a span length feature. The span length feature is trivial to compute in a constant time. However, computing the weighted sum of each span can be time consuming if not carefully designed. We decompose the computation as a set of summation-based operations (see Algorithm 2 and 3) so as to use summed area tables (Szeliski, 2010), which was been used in Li et al. (2019) for constant time computation of span representations. These pseudocode definitions are designed based on Tensor operations, which are highly optimized and fast.

---

[6]https : / / github . com / tensorflow / tensor2tensor/blob/master/tensor2tensor/ layers/area_attention.py

---

**Algorithm 2:** Compute the weighted embedding sum of each span in parallel, using ComputeSpanVectorSum defined in Algorithm 3.

---

**Input:** Tensors $H$ and $E$ are the hidden and embedding vectors of a sequence of tokens respectively, in shape of $[L, D]$ with length $L$ and depth $D$.

**Output:** weighted embedding sum, $\hat{X}$.

1 **Hyperparameter**: max span length $M$.
2 Compute token weights $A$:
    $A \leftarrow \exp(\phi(H, \theta)W)$ where $\phi(\cdot)$ is a multi-layer perceptron with trainable parameters $\theta$, followed by a linear transformation $W$. $A \in R^{L \times 1}$;
3 $E' \leftarrow E \otimes A$ where $\otimes$ is element-wise multiplication. The last dim of $A$ is broadcast;
4 $\hat{E} \leftarrow$ ComputeSpanVectorSum$(E')$;
5 $\hat{A} \leftarrow$ ComputeSpanVectorSum$(A)$;
6 $\hat{X} \leftarrow \hat{E} \oslash \hat{A}$ where $\oslash$ is element-wise division. The last dim of $\hat{A}$ is broadcast;
7 **return** $\hat{X}$.

---

**Algorithm 3:** ComputeSpanVectorSum.

---

**Input:** A tensor $G$ in shape of $[L, D]$.

**Output:** Sum of vectors of each span, $U$.

1 **Hyperparameter**: max span length $M$.
2 **Compute** integral image $I$ by cumulative sum along the first dimension over $G$;
3 $I \leftarrow [0 \ I]$, padding zero to the left;
4 **for** $m = 0, \cdots, M - 1$ **do**
5     $I_1 \leftarrow I[m + 1 :, :]$ ;
6     $I_2 \leftarrow I[: -m - 1, :]$ ;
7     $\bar{I} \leftarrow I_1 - I_2$ ;
8     $U \leftarrow [U \ \bar{I}]$, concat on the first dim;
9 **return** $U$.

## C Details for Distance GCN

Given the structural distance between two objects, based on the view hierarchy tree, we compute the strength of how these objects should affect each other by applying a Gaussian kernel to the distance, as shown the following (Equation 13).

$$\text{Adjacency}(o_i, o_j) = \frac{1}{\sqrt{2\pi\sigma^2}}\exp(-\frac{d(o_i, o_j)^2}{2\sigma^2})$$

(13)

where $d(o_i, o_j)$ is the distance between object $o_i$ and $o_j$, and $\sigma$ is a constant. With this definition of soft adjacency, the rest of the computation follows the typical GCN (Niepert et al., 2016).

## D Hyperparameters & Training

We tuned all the models on a number of hyperparameters, including the token embedding depth, the hidden size and the number of hidden layers, the learning rate and schedule, and the dropout ratios. We ended up using 128 for the embedding and hidden size for all the models. Adding more dimensions does not seem to improve accuracy and slows down training.

For the phrase tuple extraction task, we used 6 hidden layers for Transformer encoder and decoder, with 8-head self and encoder-decoder attention, for all the model configurations. We used 10% dropout ratio for attention, layer preprocessing and relu dropout in Transformer. We followed the learning rate schedule detailed previously (Vaswani et al., 2017), with an increasing learning rate to 0.001 for the first 8K steps followed by an exponential decay. All the models were trained for 1 million steps with a batch size of 128 on a single Tesla V100 GPU, which took 28 to 30 hours.

For the grounding task, Filter-1 GCN and Distance GCN used 6 hidden layers with ReLU for nonlinear activation and 10% dropout ratio at each layer. Both GCN models use a smaller peak learning rate of 0.0003. The Transformer screen encoder also uses 6 hidden layers but uses a much larger dropout ratio: ReLU dropout of 30%, attention dropout of 40%, and layer preprocessing dropout of 20%, with a peak learning rate of 0.001. All the grounding models were trained for 250K steps on the same hardware.