

## (D) Pay Attention (1/2) [10 Points]

The meaning of a word depends on its context. For example, in the sentence “The farmer *seeded* the field with corn,” the word *seeded* means “added seeds to.” However, in the sentence “The chef *seeded* the tomato,” the word *seeded* means “took seeds away from.” Similarly, *bat* means different things in “The baseball player swung the *bat*” vs. “The *bat* flew through the air.”

If you were building a model of language, how would you get it to recognize the way that a word’s meaning depends on context? One popular technique for achieving this goal is a mechanism called **attention**. In the way that attention is implemented in current state-of-the-art models of language, the model has a large number of **attention heads**, each of which is denoted by a pair of numbers (for examples, 8-10). When the model processes a sentence, for every pair of words in the sentence, each head calculates the “relatedness” of the two words.

The one wrinkle is that we do not know what exactly “relatedness” should mean, so instead of telling the model how to define “relatedness,” we let the model learn its own definition of relatedness. Recently, computer scientists have started to analyze what these attention heads have learned, and this analysis shows that they often reflect linguistic information! For example, here’s the output of one attention head (head 8-10) when we feed the following sentence into BERT<sup>1</sup>, which is one of the most popular models that uses attention heads:

**Example:** I<sub>1</sub> see<sub>2</sub> my<sub>3</sub> sister<sub>4</sub>, but<sub>5</sub> she<sub>6</sub> can’t<sub>7</sub> see<sub>8</sub> me<sub>9</sub> because<sub>10</sub> she<sub>11</sub> is<sub>12</sub> reading<sub>13</sub> a<sub>14</sub> linguistics<sub>15</sub> book<sub>16</sub>.

**Output:** 4 → 2, 9 → 8, 16 → 13

This output signifies that head 8-10 connects word #4 (*sister*) to word #2 (*see*), as well as word #9 (*me*) to word #8 (*see*) and word #16 (*book*) to word #13 (*reading*). If you consider what all of those pairs of words have in common, you’ll see that each one is a verb and its direct object: *sister* is the direct object of the first instance of *see*, *me* is the direct object of the second instance of *see*, and *book* is the direct object of *reading*. It appears that head 8-10 has learned to connect verbs to their objects! (Note that these activations are directional; for example, word #2 is not connected to word #4.)

Why is this information useful? If we go back to the example with the verb *seed*, this sort of information can help the model figure out which version of *seed* is being used: If its direct object is something like *field* or *lawn*, then it probably means “to add seeds to;” if its direct object is something like *tomato* or *watermelon*, then it probably means “to take seeds away from.” Of course, one sentence isn’t enough to draw strong conclusions. Instead, computer scientists tend to use a *corpus*, or a database, of example sentences to find patterns in the data. On the next page is a small corpus, the NacloWeb Corpus<sup>1</sup>, which has 7 sentences.

<sup>1</sup> Some sentences derived from data in the English Web Treebank.



## (D) Pay Attention (2/2)

### NacloWeb Corpus

1. My<sub>1</sub> experience<sub>2</sub> with<sub>3</sub> Gelda<sub>4</sub> 's<sub>5</sub> House<sub>6</sub> of<sub>7</sub> Gelbelgarg<sub>8</sub> has<sub>9</sub> been<sub>10</sub> extremely<sub>11</sub> wonderful<sub>12</sub>
2. We<sub>1</sub> use<sub>2</sub> Google<sub>3</sub> 's<sub>4</sub> models<sub>5</sub> to<sub>6</sub> delve<sub>7</sub> into<sub>8</sub> the<sub>9</sub> inner<sub>10</sub> workings<sub>11</sub> of<sub>12</sub> language<sub>13</sub>
3. At<sub>1</sub> this<sub>2</sub> corporation<sub>3</sub> 's<sub>4</sub> meeting<sub>5</sub> people<sub>6</sub> are<sub>7</sub> concerned<sub>8</sub> about<sub>9</sub> the<sub>10</sub> company<sub>11</sub> 's<sub>12</sub> plans<sub>13</sub>
4. In<sub>1</sub> July<sub>2</sub> we<sub>3</sub> will<sub>4</sub> interview<sub>5</sub> the<sub>6</sub> candidate<sub>7</sub> and<sub>8</sub> review<sub>9</sub> her<sub>10</sub> resumé<sub>11</sub> again<sub>12</sub>
5. The<sub>1</sub> platypus<sub>2</sub> is<sub>3</sub> a<sub>4</sub> strange<sub>5</sub> animal<sub>6</sub>, with<sub>7</sub> its<sub>8</sub> eggs<sub>9</sub> and<sub>10</sub> its<sub>11</sub> webbed<sub>12</sub> feet<sub>13</sub>
6. I<sub>1</sub> think<sub>2</sub> that<sub>3</sub> although<sub>4</sub> my<sub>5</sub> NACLO<sub>6</sub> exam<sub>7</sub> was<sub>8</sub> difficult<sub>9</sub>, it<sub>10</sub> was<sub>11</sub> a<sub>12</sub> lot<sub>13</sub> of<sub>14</sub> fun<sub>15</sub>
7. Linguistics<sub>1</sub> is<sub>2</sub> a<sub>3</sub> beautiful<sub>4</sub> science<sub>5</sub> that<sub>6</sub> provides<sub>7</sub> interdisciplinary<sub>8</sub> insight<sub>9</sub> into<sub>10</sub> the<sub>11</sub> human<sub>12</sub> experience<sub>13</sub>

Note that the NacloWeb corpus treats the possessive element 's as a separate word. (So in Sentence #1, word #5 is 's and word #6 is House.)

In our experiment on the NacloWeb Corpus, we ran each of the corpus' sentences through BERT and recorded the outputs of four attention heads (8-11, 7-6, 9-6, and 5-4). Unfortunately, due to some extremely sloppy experimental procedure, we don't remember in which order we ran them through the model; in addition, we forgot to record some data. Your job is to fill in the blanks! Note that some blanks may have more than one connection, and some may have none at all.

Sentence	8-11	7-6	9-6	5-4
Sentence A	12 → 13	5 → 7	14 → 15	10 → 7, 5 → 1
Sentence B	(a)	8 → 9, 11 → 13	7 → 9, 7 → 13	8 → 2, 11 → 2
Sentence C	2 → 3, 10 → 11	4 → 5, 12 → 13	(b)	11 → 3
Sentence D	(c)	4 → 5	12 → 13, 8 → 11	None
Sentence E	3 → 5, 11 → 13	(d)	10 → 13	5 → 1
Sentence F	(e)	1 → 2, 5 → 6	7 → 8, 3 → 6	None
Sentence G	(f)	(g)	(h)	(i)

**D1.** Identify sentences A-G. Record your answers in the Answer Sheets.

**D2.** Fill in the missing data in the table in the Answer Sheets.

