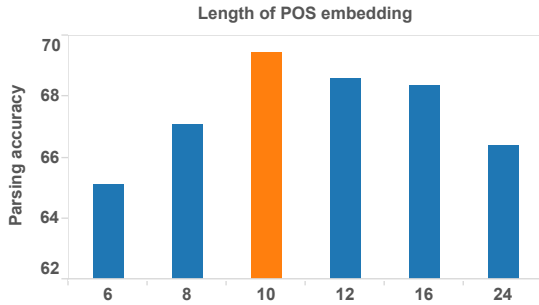


	Arabic	Basque	Czech	Danish	Dutch	Portuguese	Slovene	Swedish
Standard EM								
DMV	<b>45.8</b>	41.1	31.3	50.8	47.1	<b>36.7</b>	36.7	43.5
Neural DMV	43.4	<b>46.5</b>	<b>33.1</b>	<b>55.6</b>	<b>49.0</b>	30.4	<b>42.2</b>	<b>44.3</b>
Softmax EM $\sigma = 0.25$								
DMV	49.3	45.6	30.4	43.6	<b>46.1</b>	33.5	29.8	<b>50.3</b>
Neural DMV	<b>54.2</b>	<b>46.3</b>	<b>36.8</b>	<b>44.0</b>	39.9	<b>35.8</b>	<b>31.2</b>	49.7
Softmax EM $\sigma = 0.5$								
DMV	<b>54.2</b>	47.6	<b>43.2</b>	38.8	<b>38.0</b>	33.7	23.0	37.2
Neural DMV	44.6	<b>48.9</b>	33.4	<b>50.3</b>	37.5	<b>35.3</b>	<b>32.2</b>	<b>43.3</b>
Softmax EM $\sigma = 0.75$								
DMV	42.2	<b>48.6</b>	22.7	41.0	<b>33.8</b>	33.5	<b>23.2</b>	41.6
Neural DMV	<b>56.7</b>	45.3	<b>31.6</b>	<b>41.3</b>	33.7	<b>34.7</b>	22.9	<b>42.0</b>
Viterbi EM								
DMV	32.5	47.1	27.1	39.1	37.1	32.3	23.7	42.6
Neural DMV	<b>48.2</b>	<b>48.1</b>	<b>28.6</b>	<b>39.8</b>	<b>37.2</b>	<b>36.5</b>	<b>39.9</b>	<b>47.9</b>

**Table 3:** DDA results (on sentences no longer than 10) on eight additional languages. Our neural based approaches are compared with traditional approaches using standard EM, softmax EM (parameterized by  $\sigma$ ) and Viterbi EM.

Activation function	WSJ10
ReLU	<b>69.7</b>
Leaky ReLU	67.0
Tanh	66.2
Sigmoid	62.5
Linear	55.1

**Table 5:** Comparison between activation functions.



**Figure 3:** Parsing accuracy vs. length of POS embedding

different POS tags. On the other hand, when the dimension is too high (such as  $dim = 30$ ), since we have only 35 POS tags, the neural network is prone to overfitting.

**Shared parameters** An alternative to our neural network architecture is to have two separate neural networks to compute *CHILD* and *DECISION* rule probabilities respectively. The embeddings of the head POS tag and the valence are not shared between the two networks. As can be seen in Table

	WSJ10	WSJ
Separate Networks	68.6	52.1
Merged Network	<b>69.7</b>	<b>52.5</b>

**Table 6:** Comparison between using two separate networks and using a merged network.

6, sharing POS tags embeddings attribute to better performance.

## 5 Model Analysis

In this section, we investigate what information our neural based DMV model captures and analyze how it contributes to better parsing performance.

### 5.1 Correlation of POS Tags Encoded in Embeddings

A main motivation of our approach is to encode correlation between POS tags in their embeddings so as to smooth the probabilities of grammar rules involving correlated POS tags. Here we want to examine whether the POS embeddings learned by our approach successfully capture such correlation.

We collected the POS embeddings learned in the experiment described in section 4.3 and visualized them on a 2D plane using the t-SNE algorithm (Van der Maaten and Hinton, 2008). t-SNE is a dimensionality reduction algorithm that maps data from a high dimensional space to a low dimensional one (2 or 3) while maintaining the distances between