

# Joint Lemmatization and Morphological Tagging with LEMMING Appendix

Section 1 and 2 illustrate the edit tree extraction and application by Chrupała (2008) and also provide pseudo code. Section 3 and 4 provide an extended overview of our results.

## 1 Tree Extraction

```

1: function TREE( $x, y$ )
2:    $i_s, i_e, j_s, j_e \leftarrow LCS(x, y)$ 
3:   if  $i_e - i_s = 0$  then
4:     return SUB( $x, y$ )
5:   else
6:     return (TREE( $x_0^{i_s}, y_0^{j_s}$ ),  $i_s$ , TREE( $x_{i_e}^{|x|}, y_{j_e}^{|y|}$ ),  $|x| - i_e$ )

```

Create a tree given a form-lemma pair  $\langle x, y \rangle$ . LCS returns the start and end indexes of the LCS in  $x$  and  $y$ .  $x_{i_s}^{i_e}$  denotes the substring of  $x$  starting at index  $i_s$  (inclusive) and ending at index  $i_e$  (exclusive).  $i_e - i_s$  thus equals the length of this substring.  $|x|$  denotes the length of  $x$ . Note that the tree does not store the LCS, but only the length of the prefix and suffix. This way the tree for *umgeschaut* can also be applied to transform *umgebaut* “renovated” into *umbauen* “to renovate”.

For the example *umgeschaut-umschauen*, the LCS is the stem *schau*. The function then recursively transforms *umge* into *um* and *t* into *en*. The prefix and suffix lengths of the form are 4 and 1 respectively. The left sub-node needs to transform *umge* into *um*. The new LCS is *um*. The new prefix and suffix lengths are 0 and 2 respectively. As the new prefix is empty there is nothing more to do. The suffix node needs to transform *ge* into the empty string  $\epsilon$ . As the new LCS of the suffix is empty, because *ge* and  $\epsilon$  have no character in common, the node is represented as a substitution node. The remaining transformation *t* into *en* is also represented as a substitution, resulting in the tree in Figure 1:

## 2 Tree Application

```

1: function APPLY(tree,  $x$ )

```

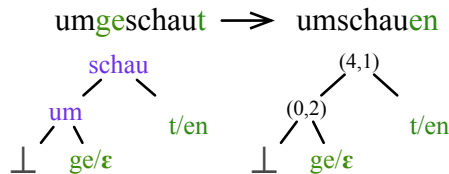


Figure 1: Edit tree for the inflected form *umgeschaut* “looked around” and its lemma *umschauen* “to look around”. The right tree is the actual edit tree we use in our model, the left tree visualizes what each node corresponds to. Note how the root node stores the length of the prefix *umge* and the suffix *t*.

```

2:   if tree is a LCS node then
3:     tree  $\rightarrow$  treei, il, treej, jl
4:     if |x| < il + jl then                                      $\triangleright$  Prefix and Suffix do not fit.
5:       return  $\perp$ 
6:       p = APPLY(treei, x0i)                                      $\triangleright$  Create prefix.
7:       if p is  $\perp$  then                                            $\triangleright$  Prefix tree cannot be applied.
8:         return  $\perp$ 
9:         s = APPLY(treej, x|x||x| - jl)                                      $\triangleright$  Create suffix.
10:      if s is  $\perp$  then                                            $\triangleright$  Suffix tree cannot be applied.
11:        return  $\perp$ 
12:      return p + xi|x| - jl + s                                $\triangleright$  Concatenate prefix, LCS and suffix.
13:    else                                                            $\triangleright$  tree is a SUB node
14:      tree  $\rightarrow$  u, v
15:      if x = u then                                            $\triangleright$  If x and u match return v
16:        return v
17:      return  $\perp$                                                   $\triangleright$  tree cannot be applied.

```

In the code + represents string concatenation and  $\perp$  a null string, meaning that the tree cannot be applied to the form. We first run the tree depicted in Figure 1 on the form *angebaut* “attached (to a building)”. The first node is a LCS node specifying that prefix and suffix should have length 4 and 1, respectively. We thus recursively apply the left child node to the prefix string *ange*. This is done by matching the length two prefix *an* and deleting *ge* yielding the intermediate result *anbaut*. We continue on the right side of the tree and replace *t* with *en*. This yield the final (correct) result *anbauen*. The application of the tree to the form *einbauen* “installed” would fail, as we would try to substitute a *ge* in *eing*.

### 3 Development Results

		cs	de	en	es	hu	la							
Baselines	SIMPLE tag	86.08	69.83	79.05	56.02	94.72	87.95	96.29	87.21	94.37	85.20	83.86	57.77	
	MORFETTE lemma	87.22	32.95	93.27	62.15	97.60	75.64	92.92	59.68	86.09	42.02	85.19	14.06	
	joint	77.18	22.56	75.60	36.43	93.15	67.58	90.28	54.00	82.99	37.03	75.00	5.39	
	SIMPLE tag	<u>89.82</u>	<u>76.83</u>	<u>85.05</u>	<u>65.22</u>	<b>95.71</b>	<b>90.29</b>	<u>96.83</u>	<u>89.00</u>	<b>95.46</b>	<u>88.17</u>	<u>86.35</u>	<u>65.21</u>	
	PCRF lemma	87.36	32.95	93.28	62.15	97.66	75.64	92.99	59.68	86.11	42.02	85.35	14.06	
	joint	79.99	23.98	80.71	40.26	94.09	69.48	90.75	54.66	83.47	37.16	76.58	6.61	
	JCK tag	86.08	69.83	79.05	56.02	94.72	87.95	96.29	87.21	94.37	85.20	83.86	57.77	
	MORFETTE lemma	96.24	82.59	97.67	88.80	98.71	92.50	97.61	86.76	97.48	91.16	<u>93.26</u>	<u>63.09</u>	
	joint	84.32	61.67	78.10	51.42	94.43	84.61	94.67	78.39	93.15	80.73	81.34	43.71	
	JCK tag	<u>89.82</u>	<u>76.83</u>	<u>85.05</u>	<u>65.22</u>	<b>95.71</b>	<b>90.29</b>	<u>96.83</u>	<u>89.00</u>	<b>95.46</b>	<u>88.17</u>	<u>86.35</u>	<u>65.21</u>	
	PCRF lemma	<u>96.26</u>	<u>82.88</u>	<u>97.74</u>	89.11	<u>98.81</u>	<u>93.06</u>	97.78	87.10	<u>97.68</u>	<u>91.91</u>	93.06	62.64	
	joint	<u>88.17</u>	<u>68.43</u>	<u>84.16</u>	<u>60.39</u>	<u>95.41</u>	<u>86.95</u>	<u>95.37</u>	<u>80.25</u>	<u>94.33</u>	<u>83.70</u>	<u>83.57</u>	<u>48.84</u>	
	MORFETTE tag	86.08	69.83	79.05	56.02	94.72	87.95	96.29	87.21	94.37	85.20	83.86	57.77	
	MORFETTE lemma	96.25	82.54	97.12	<u>89.90</u>	98.43	92.54	<u>97.97</u>	<u>89.94</u>	97.22	90.04	91.89	55.13	
	joint	84.39	61.94	77.60	52.87	94.16	84.70	<u>95.09</u>	<u>81.95</u>	93.11	80.60	80.09	36.39	
	LEMMING	edit tree tag	86.08	69.83	79.05	56.02	94.72	87.95	96.29	87.21	94.37	85.20	83.86	57.77
		MORFETTE lemma	96.29	82.93	97.84	89.78	98.71	92.63	97.91	89.00	97.31	90.49	93.00	61.68
		joint	84.45	62.36	78.32	52.75	94.43	84.79	94.97	80.68	93.08	80.48	80.92	41.27
align tag		86.08	69.83	79.05	56.02	94.72	87.95	96.29	87.21	94.37	85.20	83.86	57.77	
MORFETTE lemma		96.74	85.38	98.17	91.61	98.76	93.11	98.05	90.13	97.70	92.15	93.76	66.30	
joint		84.72	63.88	78.47	53.52	94.47	85.09	95.05	81.29	93.31	81.51	81.49	44.67	
dict tag		86.08	69.83	79.05	56.02	94.72	87.95	96.29	87.21	94.37	85.20	83.86	57.77	
MORFETTE lemma		97.50	89.38	98.36	92.66	98.84	94.02	98.39	92.56	97.98	93.30	94.64	71.57	
joint		85.06	65.60	78.53	53.91	94.53	85.83	95.27	82.89	93.46	82.10	81.87	46.92	
morph tag		86.08	69.83	79.05	56.02	NA	NA	96.29	87.21	94.37	85.20	83.86	57.77	
MORFETTE lemma		96.59	87.28	97.43	89.96	NA	NA	98.46	92.98	97.77	92.62	93.60	66.30	
joint		85.64	67.73	78.87	55.17	NA	NA	95.63	84.69	94.10	84.02	82.29	48.72	
edit tree tag		89.82	76.83	85.05	65.22	<b>95.71</b>	<b>90.29</b>	96.83	89.00	<b>95.46</b>	<u>88.17</u>	<u>86.35</u>	<u>65.21</u>	
PCRF lemma		96.33	83.04	97.93	90.05	98.82	93.02	98.14	89.61	97.53	91.22	92.85	60.98	
joint		88.23	68.85	84.41	61.89	95.40	87.00	95.75	83.06	94.30	83.58	83.20	46.73	
align tag		89.82	76.83	85.05	65.22	<b>95.71</b>	<b>90.29</b>	96.83	89.00	<b>95.46</b>	<u>88.17</u>	<u>86.35</u>	<u>65.21</u>	
PCRF lemma		96.80	85.72	98.24	91.94	98.88	93.71	98.27	90.75	97.91	92.91	93.50	64.89	
joint		88.57	70.69	84.59	62.95	95.45	87.47	95.81	83.64	94.52	84.60	83.83	50.32	
dict tag		89.82	76.83	85.05	65.22	<b>95.71</b>	<b>90.29</b>	96.83	89.00	<b>95.46</b>	<u>88.17</u>	<u>86.35</u>	<u>65.21</u>	
PCRF lemma		97.62	89.88	98.43	92.95	98.94	<b>94.28</b>	98.62	93.23	98.23	94.16	94.49	70.60	
joint		88.97	72.79	84.63	63.16	<b>95.50</b>	<b>88.08</b>	96.02	85.24	94.73	85.46	84.41	53.79	
morph tag		89.82	76.83	85.05	65.22	NA	NA	96.83	89.00	<b>95.46</b>	<u>88.17</u>	<u>86.35</u>	<u>65.21</u>	
PCRF lemma		97.38	89.55	98.24	92.42	NA	NA	<b>98.71</b>	93.97	98.30	94.38	94.39	70.54	
joint		89.30	74.32	84.81	64.08	NA	NA	96.17	86.29	95.15	86.87	84.62	55.07	
dict tag		<b>90.47</b>	78.37	85.31	65.94	95.66	88.60	96.94	89.25	95.29	87.55	<b>86.77</b>	65.98	
joint lemma		98.34	92.96	98.66	94.07	<b>99.00</b>	94.23	98.68	94.08	98.53	95.46	<b>96.68</b>	<b>82.99</b>	
joint		89.82	75.47	84.94	64.19	95.48	86.52	96.11	85.43	94.71	85.40	85.85	60.98	
morph tag		90.35	<b>79.61</b>	<b>85.52</b>	<b>67.69</b>	95.66	88.60	<b>96.95</b>	<b>89.55</b>	<b>95.46</b>	<b>88.85</b>	86.70	<b>66.82</b>	
joint lemma		<b>98.55</b>	<b>94.12</b>	<b>98.74</b>	<b>94.47</b>	<b>99.00</b>	94.23	98.68	<b>94.53</b>	<b>98.57</b>	<b>95.65</b>	96.51	82.22	
joint		<b>90.05</b>	<b>78.37</b>	<b>85.28</b>	<b>66.58</b>	95.48	86.52	<b>96.20</b>	<b>86.51</b>	<b>95.30</b>	<b>88.22</b>	<b>85.97</b>	<b>62.97</b>	

Table A1: Development accuracies for the baselines and the different pipeline versions and a joint version of LEMMING. The numbers in each cell are general token accuracy and the token accuracy on unknown forms. Each cell specifies either a baseline  $\in$  {baseline, JCK, MORFETTE} or a LEMMING feature set  $\in$  {edit tree, align, dict, morph} and a tagger  $\in$  {MORFETTE, PCRF, joint}.

## 4 Test Results

		cs	de	en	es	hu	la							
Baselines	SIMPLE tag	86.00	68.88	76.80	52.86	95.35	87.58	96.30	87.65	92.33	81.32	79.70	47.06	
	MORFETTE lemma	86.81	31.16	91.42	58.68	97.88	78.95	92.86	59.34	84.67	38.95	83.60	20.43	
	MORFETTE joint	76.96	21.22	72.83	33.27	93.82	68.98	90.14	53.45	80.41	33.08	71.62	7.52	
	SIMPLE tag	<u>89.75</u>	<u>76.83</u>	<u>82.81</u>	<u>61.60</u>	<b>96.45</b>	<b>90.68</b>	<u>97.05</u>	<u>90.07</u>	<u>93.64</u>	<u>84.65</u>	<u>82.37</u>	<u>53.73</u>	
	PCRF lemma	86.94	31.16	91.48	58.68	97.92	78.95	92.92	59.34	84.73	38.95	83.80	20.43	
	PCRF joint	79.66	22.86	77.78	36.49	94.87	71.62	90.81	54.98	81.06	33.27	73.51	8.58	
	JCK tag	86.00	68.88	76.80	52.86	95.35	87.58	96.30	87.65	92.33	81.32	79.70	47.06	
	MORFETTE lemma	95.87	80.79	96.50	85.54	98.95	93.76	97.58	86.80	96.52	88.09	<u>90.84</u>	58.13	
	MORFETTE joint	84.00	59.46	75.60	47.77	95.09	84.34	94.54	77.97	90.71	75.62	76.76	33.56	
	JCK tag	<u>89.75</u>	<u>76.83</u>	<u>82.81</u>	<u>61.60</u>	<b>96.45</b>	<b>90.68</b>	<u>97.05</u>	<u>90.07</u>	<u>93.64</u>	<u>84.65</u>	<u>82.37</u>	<u>53.73</u>	
	PCRF lemma	<u>95.95</u>	<u>81.28</u>	<u>96.63</u>	<u>85.84</u>	<u>99.08</u>	<u>94.28</u>	<u>97.69</u>	<u>87.19</u>	<u>96.69</u>	<u>88.66</u>	<u>90.79</u>	<u>58.23</u>	
	PCRF joint	<u>87.85</u>	<u>67.00</u>	<u>81.60</u>	<u>55.97</u>	<u>96.17</u>	<u>87.32</u>	<u>95.44</u>	<u>80.62</u>	<u>92.15</u>	<u>78.89</u>	<u>79.51</u>	<u>39.07</u>	
	MORFETTE tag	86.00	68.88	76.80	52.86	95.35	87.58	96.30	87.65	92.33	81.32	79.70	47.06	
	MORFETTE lemma	95.88	80.70	95.88	<u>87.55</u>	98.73	93.56	<u>98.02</u>	<u>90.16</u>	96.06	86.21	90.08	54.16	
	MORFETTE joint	84.10	59.92	74.87	49.22	94.89	84.44	95.03	<u>81.68</u>	90.60	75.11	76.57	32.50	
	LEMMING	edit tree tag	86.00	68.88	76.80	52.86	95.35	87.58	96.30	87.65	92.33	81.32	79.70	47.06
		MORFETTE lemma	95.91	81.07	96.82	87.31	98.95	93.69	97.87	88.97	96.15	86.68	91.20	59.87
		MORFETTE joint	84.15	60.27	75.86	49.22	95.08	84.37	94.86	80.38	90.50	74.86	77.00	34.73
align tag		86.00	68.88	76.80	52.86	95.35	87.58	96.30	87.65	92.33	81.32	79.70	47.06	
MORFETTE lemma		96.43	83.83	97.22	89.25	98.99	94.15	98.04	90.28	96.87	89.72	91.64	62.36	
MORFETTE joint		84.47	61.97	76.05	50.07	95.11	84.73	94.98	81.30	90.94	76.76	77.18	35.79	
dict tag		86.00	68.88	76.80	52.86	95.35	87.58	96.30	87.65	92.33	81.32	79.70	47.06	
MORFETTE lemma		97.24	88.13	97.58	91.03	99.07	95.06	98.32	92.36	97.28	91.32	92.97	69.14	
MORFETTE joint		84.88	64.14	76.13	50.48	95.19	85.55	95.19	82.91	91.14	77.48	77.52	37.59	
morph tag		86.00	68.88	76.80	52.86	NA	NA	96.30	87.65	92.33	81.32	79.70	47.06	
MORFETTE lemma		96.52	86.69	96.57	88.26	NA	NA	98.53	93.73	96.70	89.30	91.53	62.31	
MORFETTE joint		85.57	66.73	76.56	51.90	NA	NA	95.63	85.02	91.92	79.73	78.18	39.65	
edit tree tag		<u>89.75</u>	<u>76.83</u>	<u>82.81</u>	<u>61.60</u>	<b>96.45</b>	<b>90.68</b>	<u>97.05</u>	<u>90.07</u>	<u>93.64</u>	<u>84.65</u>	<u>82.37</u>	<u>53.73</u>	
PCRF lemma		96.05	81.68	96.99	87.69	99.08	94.25	98.00	89.56	96.38	87.54	91.37	60.88	
PCRF joint		88.03	68.04	81.96	57.91	96.17	87.38	95.80	83.42	91.93	78.15	79.75	40.34	
align tag		<u>89.75</u>	<u>76.83</u>	<u>82.81</u>	<u>61.60</u>	<b>96.45</b>	<b>90.68</b>	<u>97.05</u>	<u>90.07</u>	<u>93.64</u>	<u>84.65</u>	<u>82.37</u>	<u>53.73</u>	
PCRF lemma		96.56	84.49	97.38	89.68	99.12	94.51	98.17	90.89	97.11	90.61	91.80	63.26	
PCRF joint		88.41	70.03	82.20	59.02	96.19	87.61	95.94	84.38	92.41	80.14	80.03	41.87	
dict tag		<u>89.75</u>	<u>76.83</u>	<u>82.81</u>	<u>61.60</u>	<b>96.45</b>	<b>90.68</b>	<u>97.05</u>	<u>90.07</u>	<u>93.64</u>	<u>84.65</u>	<u>82.37</u>	<u>53.73</u>	
PCRF lemma		97.46	89.14	97.70	91.27	<b>99.21</b>	<b>95.59</b>	98.48	92.98	97.53	92.10	93.07	69.83	
PCRF joint		88.86	72.51	82.27	59.42	<b>96.27</b>	<b>88.49</b>	96.12	85.80	92.59	80.77	80.49	44.26	
morph tag		<u>89.75</u>	<u>76.83</u>	<u>82.81</u>	<u>61.60</u>	NA	NA	97.05	90.07	93.64	84.65	82.37	53.73	
PCRF lemma		97.29	88.98	97.51	90.85	NA	NA	98.68	94.32	97.53	92.15	92.54	67.81	
PCRF joint		89.23	74.24	82.49	60.42	NA	NA	96.35	87.25	93.11	82.56	80.67	45.21	
dict tag		<b>90.34</b>	78.47	<b>83.10</b>	62.36	96.32	89.70	97.11	90.13	93.64	84.78	82.89	54.69	
joint lemma		98.27	92.67	<b>98.10</b>	92.79	<b>99.21</b>	95.23	98.67	94.07	98.02	94.15	<b>95.58</b>	<b>81.74</b>	
joint joint		89.69	75.44	82.64	60.49	96.17	87.87	96.23	86.19	92.84	81.89	81.92	49.97	
morph tag		90.20	<b>79.72</b>	<b>83.10</b>	<b>63.10</b>	96.32	89.70	<b>97.17</b>	<b>90.66</b>	<b>93.67</b>	<b>85.12</b>	<b>83.49</b>	<b>58.76</b>	
joint lemma		<b>98.42</b>	<b>93.46</b>	<b>98.10</b>	<b>93.02</b>	<b>99.21</b>	95.23	<b>98.78</b>	<b>94.86</b>	<b>98.08</b>	<b>94.26</b>	95.36	80.94	
joint joint		<b>89.90</b>	<b>78.34</b>	<b>82.84</b>	<b>62.10</b>	96.17	87.87	<b>96.41</b>	<b>87.47</b>	<b>93.40</b>	<b>84.15</b>	<b>82.57</b>	<b>54.63</b>	

Table A2: Test accuracies for the baselines and the different pipeline versions and a joint version of LEMMING. The numbers in each cell are general token accuracy and the token accuracy on unknown forms. Each cell specifies either a baseline  $\in$  {baseline, JCK, MORFETTE} or a LEMMING feature set  $\in$  {edit tree, align, dict, morph} and a tagger  $\in$  {MORFETTE, PCRF, joint}.

	cs	de	es	hu	la
+dict	98.35	99.04	98.92	98.83	96.14
+morph	<b>99.03</b>	<b>99.47</b>	<b>99.09</b>	<b>99.41</b>	<b>96.73</b>

Table A3: Development accuracies for LEMMING with and without morphological attributes using *gold* tags.

	train				dev				test			
	sent	token	pos	morph	sent	token	form unk	lemma unk	sent	token	form unk	lemma unk
cs	5979	100012	12	266	5228	87988	19.86	9.79	4213	70348	19.89	9.73
de	5662	100009	51	204	5000	76704	17.11	13.53	5000	92004	19.51	15.64
en	4028	100012	46		1336	32092	8.06	6.16	1640	39590	8.52	6.56
es	3431	100027	12	226	1655	50368	13.37	8.85	1725	50630	13.41	9.16
hu	4390	100014	22	572	1051	29989	23.51	14.36	1009	19908	24.80	14.64
la	7122	59992	23	474	890	9475	17.59	6.43	891	9922	19.82	7.56

Table A4: Dataset statistics. Showing number of sentences (sent), tokens (token), POS tags (pos), morphological tags (morph) and token-based unknown form (form unk) and lemma (lemma unk) rates.

## References

Grzegorz Chrupała. *Towards a machine-learning architecture for lexical functional grammar parsing*. PhD thesis, Dublin City University, 2008.