

# Supplementary Material: Analyzing Semantic Segmentation Using Hybrid Human-Machine CRFs

Roozbeh Mottaghi  
UCLA  
roozbehm@cs.ucla.edu

Sanja Fidler, Jian Yao, Raquel Urtasun  
TTI Chicago  
{fidler,yaojian,rurtasun}@ttic.edu

Devi Parikh  
Virginia Tech  
parikh@vt.edu

This document is supplementary to the paper [2] and contains the following items:

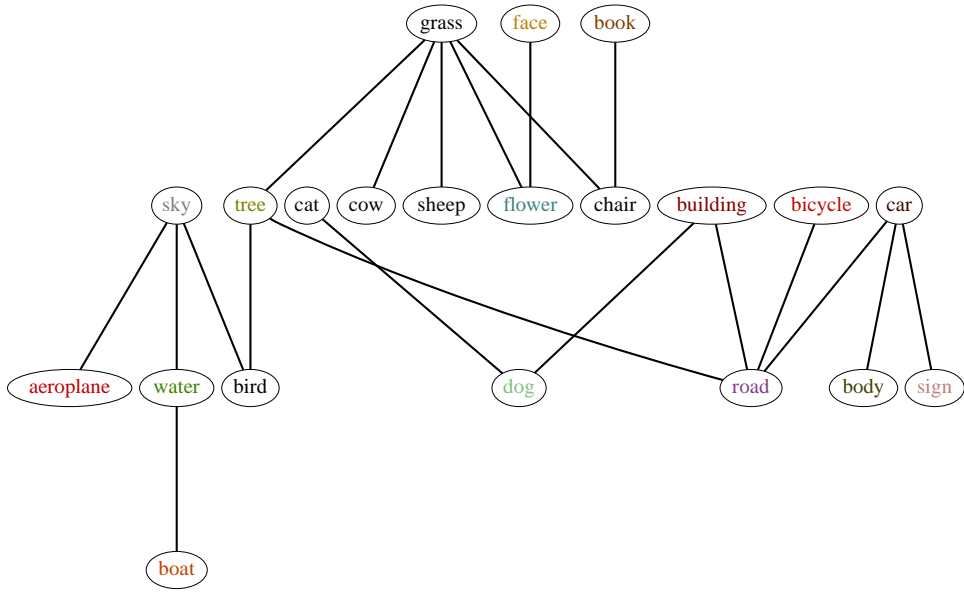
- Section 1: Human and machine co-occurrence trees
- Section 2: Per class accuracies for the plot in Figure 6
- Section 3: Human and machine confusion matrices for classification of isolated segments.

The machine model used and analyzed in the paper is [4]. Machine classification of segments in Sec.3 is obtained with Textonboost [3].

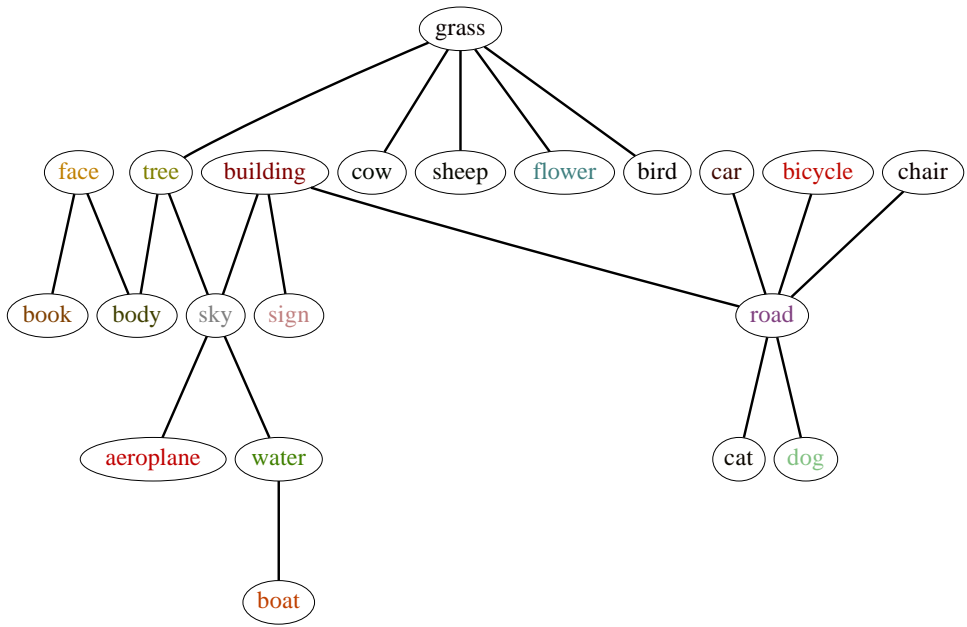
## References

- [1] C. K. Chow and C. N. Liu. Approximating discrete probability distributions with dependence trees. *IEEE Transactions on Information Theory*, 14(3):462-467, 1968. 2
- [2] R. Mottaghi, S. Fidler, J. Yao, R. Urtasun, and D. Parikh. Analyzing semantic segmentation using human-machine hybrid crfs. In *CVPR*, 2013. 1
- [3] J. Shotton, J. Winn, C. Rother, and A. Criminisi. Textonboost for image understanding: Multi-class object recognition and segmentation by jointly modeling appearance, shape and context. *IJCV*, 81(1), 2007. 1
- [4] J. Yao, S. Fidler, and R. Urtasun. Describing the scene as a whole: Joint object detection, scene classification and semantic segmentation. In *CVPR*, 2012. 1

# 1. Human and machine co-occurrence trees



(a) Human



(b) Machine

Figure 1. Chow-Liu trees [1] for humans and machine. The co-occurring categories are connected by edges. The machine tree is obtained using co-occurrence statistics from labeled training data from the MSRC dataset. The human tree is obtained by asking people (without showing them images) which pairs of categories are more likely to co-occur in a scene, and thus reflects the natural statistics of the world as perceived by humans. Both trees share several similarities.

## 2. Per class accuracies for the plot in Figure 6

In Tables 1–9, we show the per class accuracies for semantic segmentation in cases where only a single potential of the model is changed. The accuracies averaged across classes are shown in Figure 6 of the paper.

	building	grass	tree	cow	sheep	sky	aeroplane	water	face	car	bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat	average
	<b>Segments</b>																					
R	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
M	68	92	80	82	90	97	86	83	87	79	94	96	82	36	98	68	86	82	56	62	18	77.2
H	81	96	86	91	92	95	84	94	91	82	94	94	90	62	93	74	89	80	71	71	21	82.4
GT	<b>96</b>	<b>99</b>	<b>96</b>	<b>92</b>	<b>92</b>	<b>99</b>	<b>86</b>	<b>99</b>	<b>96</b>	<b>94</b>	<b>94</b>	<b>99</b>	<b>98</b>	<b>89</b>	<b>100</b>	<b>94</b>	<b>96</b>	<b>93</b>	<b>92</b>	<b>90</b>	<b>86</b>	<b>94.2</b>

Table 1. Segmentation accuracy for different types of segment potential.

	building	grass	tree	cow	sheep	sky	aeroplane	water	face	car	bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat	average
	<b>Super-segments</b>																					
R	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
M	68	92	80	82	90	97	<b>86</b>	83	87	79	94	96	82	36	98	68	86	82	56	62	18	77.2
H	<b>74</b>	<b>95</b>	83	90	90	96	84	<b>90</b>	<b>92</b>	<b>89</b>	95	97	85	<b>56</b>	99	71	<b>88</b>	81	54	71	27	81.2
GT	72	94	<b>86</b>	<b>91</b>	<b>91</b>	<b>98</b>	84	86	92	85	<b>95</b>	<b>97</b>	<b>93</b>	54	<b>99</b>	<b>83</b>	87	<b>85</b>	<b>59</b>	<b>72</b>	<b>32</b>	<b>82.6</b>

Table 2. Segmentation accuracy for different types of super-segment potential.

	building	grass	tree	cow	sheep	sky	aeroplane	water	face	car	bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat	average
	<b>Scene</b>																					
R	68	92	80	82	89	97	86	83	<b>86</b>	79	94	96	85	35	98	<b>70</b>	86	78	55	62	23	77.4
M	68	92	80	82	90	97	86	83	<b>87</b>	79	94	96	82	36	98	68	86	<b>82</b>	56	62	18	77.2
H	68	92	80	82	90	97	86	83	86	79	94	96	85	36	98	68	86	78	56	62	24	77.4
GT	<b>68</b>	<b>92</b>	<b>80</b>	<b>82</b>	<b>90</b>	<b>97</b>	<b>86</b>	<b>83</b>	86	<b>79</b>	<b>94</b>	<b>96</b>	<b>85</b>	<b>36</b>	<b>98</b>	68	<b>86</b>	79	<b>56</b>	<b>63</b>	<b>25</b>	<b>77.5</b>

Table 3. Segmentation accuracy for different types of scene potential.

	building	grass	tree	cow	sheep	sky	aeroplane	water	face	car	bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat	average
	<b>Shape</b>																					
R	68	92	80	83	90	97	86	83	84	81	94	97	83	30	98	68	86	82	56	62	10	76.6
M	68	92	80	82	90	97	86	83	87	79	94	96	82	36	98	68	86	82	56	62	18	77.2
H	69	92	80	85	91	97	85	83	84	80	94	96	82	30	98	67	87	82	55	67	22	77.4
GT	<b>69</b>	<b>92</b>	<b>80</b>	<b>91</b>	<b>91</b>	<b>97</b>	<b>86</b>	<b>83</b>	<b>93</b>	<b>84</b>	<b>94</b>	<b>97</b>	<b>83</b>	<b>36</b>	<b>98</b>	<b>73</b>	<b>86</b>	<b>82</b>	<b>57</b>	<b>61</b>	<b>49</b>	<b>80.2</b>

Table 4. Segmentation accuracy for different types of shape potential.

	building	grass	tree	cow	sheep	sky	aeroplane	water	face	car	bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat	average
	<b>Class-Class</b>																					
R	68	91	80	82	90	97	86	83	87	79	94	96	82	36	98	68	86	78	55	62	<b>19</b>	77.0
M	68	92	80	82	90	97	86	83	87	79	94	96	82	36	98	68	86	82	56	62	18	77.2
H	<b>68</b>	<b>92</b>	<b>80</b>	<b>82</b>	<b>90</b>	<b>97</b>	<b>86</b>	<b>83</b>	<b>87</b>	<b>79</b>	<b>94</b>	<b>96</b>	<b>82</b>	<b>36</b>	<b>98</b>	<b>68</b>	<b>86</b>	<b>82</b>	<b>56</b>	<b>62</b>	18	<b>77.2</b>
GT	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>

Table 5. Segmentation accuracy for different types of class-class potential.

	building	grass	tree	cow	sheep	sky	aeroplane	water	face	car	bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat	average
	<b><math>P^n</math> Potts</b>																					
R	<b>70</b>	90	<b>83</b>	<b>84</b>	<b>91</b>	<b>98</b>	85	81	<b>89</b>	77	91	93	<b>86</b>	36	98	68	<b>87</b>	<b>83</b>	56	60	10	77.0
M	68	<b>92</b>	80	82	90	97	<b>86</b>	<b>83</b>	87	<b>79</b>	<b>94</b>	<b>96</b>	82	<b>36</b>	<b>98</b>	<b>68</b>	86	82	<b>56</b>	<b>62</b>	<b>18</b>	<b>77.2</b>
H	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>
GT	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>

Table 6. Segmentation accuracy for different types of  $P^n$  Potts potential.

	building	grass	tree	cow	sheep	sky	aeroplane	water	face	car	bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat	average
	<b>Scene-Class</b>																					
R	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
M	68	92	80	82	<b>90</b>	97	86	83	<b>87</b>	79	94	96	82	<b>36</b>	98	68	86	<b>82</b>	<b>56</b>	62	18	77.2
H	<b>68</b>	<b>92</b>	<b>80</b>	<b>82</b>	89	<b>97</b>	<b>86</b>	<b>83</b>	86	<b>79</b>	<b>94</b>	<b>96</b>	<b>85</b>	35	<b>98</b>	<b>70</b>	<b>86</b>	78	55	<b>62</b>	<b>22</b>	<b>77.3</b>
GT	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X

Table 7. Segmentation accuracy for different types of scene-class potential.

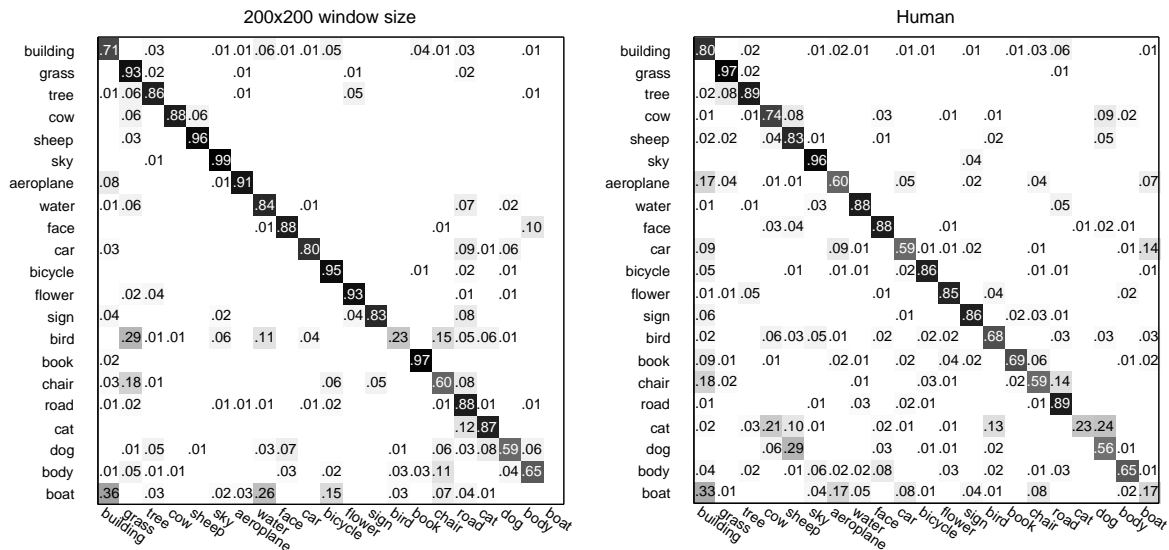
	building	grass	tree	cow	sheep	sky	aeroplane	water	face	car	bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat	average
	<b>Detector</b>																					
R	67	92	<b>81</b>	81	89	<b>98</b>	84	<b>84</b>	83	78	93	96	78	30	98	66	<b>88</b>	77	54	60	0	75.0
M	<b>68</b>	92	80	82	90	97	86	83	87	79	94	96	82	<b>36</b>	98	68	86	82	56	<b>62</b>	18	77.2
H	67	92	80	83	91	97	86	83	92	84	95	98	92	32	100	71	87	83	57	61	69	80.8
GT	67	<b>92</b>	80	<b>83</b>	<b>91</b>	97	<b>86</b>	83	<b>92</b>	<b>84</b>	<b>95</b>	<b>98</b>	<b>92</b>	32	<b>100</b>	<b>71</b>	87	<b>83</b>	<b>57</b>	61	<b>69</b>	<b>80.8</b>

Table 8. Segmentation accuracy for different types of detector potential.

	building	grass	tree	cow	sheep	sky	aeroplane	water	face	car	bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat	average
	<b>Class Presence</b>																					
R	68	92	80	84	90	97	86	84	87	80	<b>95</b>	96	83	35	98	72	86	82	56	62	18	77.6
M	68	92	80	82	90	97	86	83	87	79	94	96	82	36	98	68	86	82	56	62	18	77.2
H	68	91	81	85	<b>91</b>	97	86	83	<b>87</b>	80	94	96	86	35	98	<b>72</b>	84	<b>83</b>	55	61	18	77.8*
GT	<b>69</b>	<b>93</b>	<b>82</b>	<b>88</b>	90	<b>98</b>	<b>86</b>	<b>84</b>	84	<b>80</b>	94	<b>96</b>	<b>87</b>	<b>37</b>	<b>99</b>	70	<b>86</b>	79	<b>57</b>	<b>64</b>	<b>28</b>	<b>78.6</b>

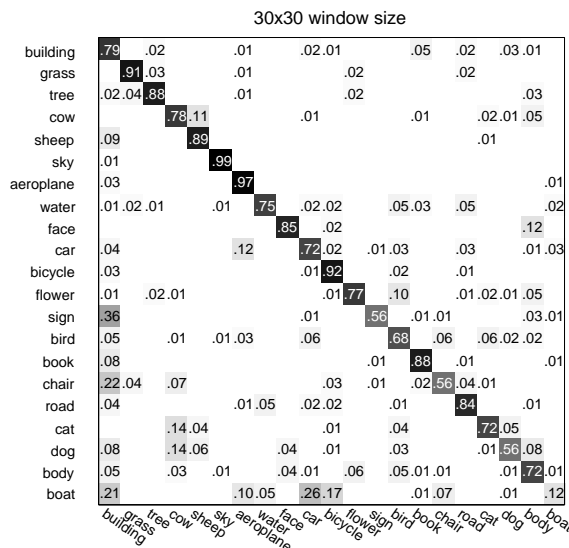
Table 9. Segmentation accuracy for different types of class presence potential. \*It should be noted this does NOT correspond to showing humans an image and asking them whether a certain object is present in the image or not. Human responses would be perfect and would correspond to GT. Instead, we show subjects a sampling of images from the dataset and ask them to build an intuition for the scenes. We then ask them which categories are likely to be more frequent than others in the dataset. We use these as potentials instead of the category occurrence frequencies extracted from training images.

### 3. Human and machine confusion matrices for classification of isolated segments



(a) Machine classification with 200x200 windows

(b) Human classification



(c) Machine classification with 30x30 windows

Figure 2. The confusion matrices for segment classification are shown. For machine with large window size (a), there is high confusion between classes appearing in the surrounding area of each other, for example, bird-grass, car-road, etc. The types of mistakes are different for humans (b). For instance, there is confusion between cat-cow or boat-aeroplane. When we reduce the window size for machine (c), the mistakes become more similar to the human mistakes. Combining the small-window machine potentials with the large-window machine potentials results in a significant improvement in segmentation accuracy, resulting in the state-of-the-art performance on the MSRC dataset.