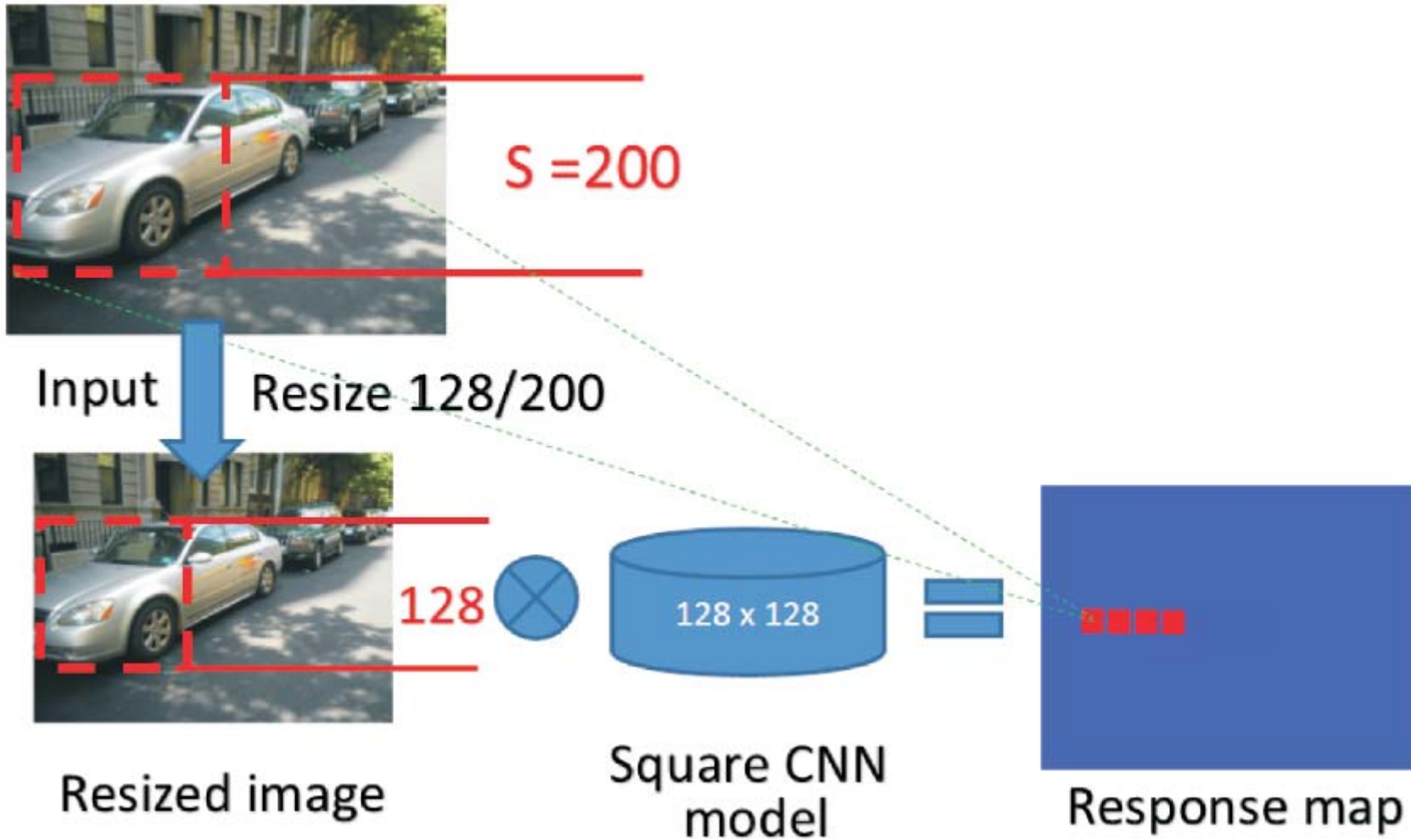


Object Image

Compact Square Object Image

Accuracy	CNN model	Human
whole image	87.3%	97.6%
compact square image	84.6%	95.4%

Table 1. Empirical recognition accuracy of the CNN model [14] and human recognition on the VOC 2007 dataset. We ask people to label the category based on compact square object (CSO) images. With more careful parameter tuning of the CNN we believe the above accuracy can be further improved.



---

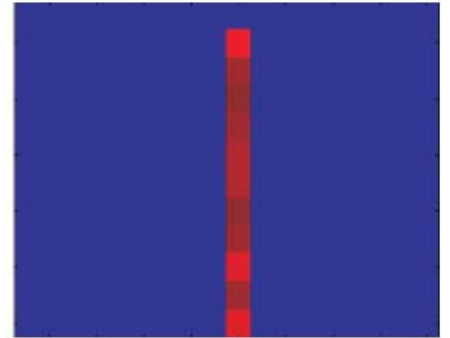
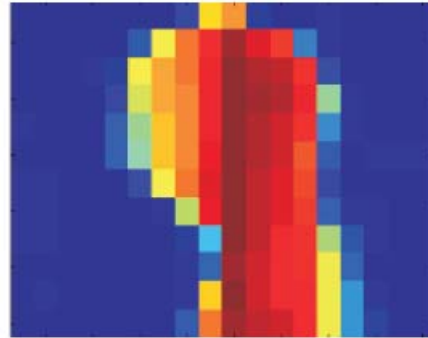
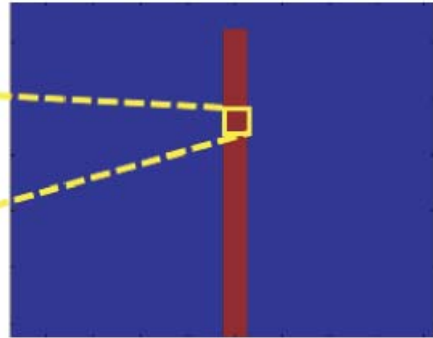
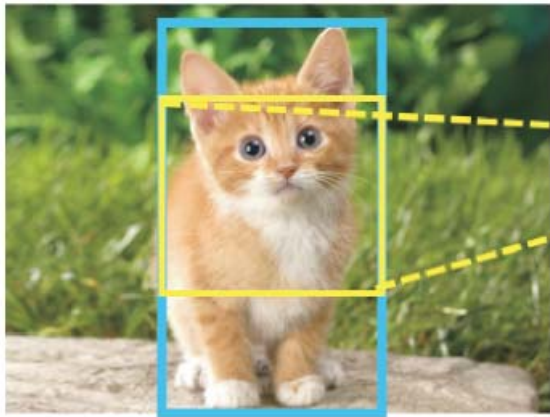
**Algorithm 1** Pseudo-code of Scale Selection Function

---

**Function** *scaleSelection*( $S, L$ )  
**if**  $d(S, L) < \varepsilon$  or  $L = S + 1$  **then**  
    **return**  
**else**  
    selecting scale  $M = \sqrt{LS}$   
    **call** *scaleSelection*( $S, M$ )  
    **call** *scaleSelection*( $M + 1, L$ )  
**end if**

---





VOC 2007 test	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mAP
R-CNN (SS)	64.2	69.7	50.0	41.9	32.0	62.6	71.0	60.7	32.7	58.5	46.5	56.1	60.6	66.8	54.2	31.5	52.8	48.9	57.9	64.7	54.2
R-CNN BB (SS)	68.1	72.8	<b>56.8</b>	43.0	36.8	66.3	74.2	67.6	34.4	63.5	54.5	61.2	69.1	68.6	58.7	33.4	62.9	51.1	62.5	64.8	58.5
R-CNN (EB)	63.7	64.8	50.3	42.2	31.4	59.4	68.4	61.3	31.2	54.7	42.6	58.6	63.1	64.2	56.3	30.2	53.1	45.2	58.6	66.2	53.2
Ours	<b>73.1</b>	<b>76.1</b>	55.8	<b>45.6</b>	<b>40.6</b>	<b>71.4</b>	<b>80.0</b>	<b>69.4</b>	<b>39.1</b>	<b>67.4</b>	<b>57.6</b>	<b>64.7</b>	<b>72.5</b>	<b>72.3</b>	<b>62.6</b>	<b>35.7</b>	<b>64.6</b>	<b>56.2</b>	<b>65.6</b>	<b>69.7</b>	<b>62.0</b>

VOC 2012 test	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mAP
R-CNN (SS)	68.1	63.8	46.1	29.4	27.9	56.6	57.0	65.9	26.5	48.7	39.5	66.2	57.3	65.4	53.2	26.2	54.5	38.1	50.6	51.6	49.6
R-CNN BB (SS)	71.8	65.8	52.0	34.1	32.6	59.6	60.0	69.8	27.6	52.0	<b>41.7</b>	69.6	61.3	68.3	57.8	29.6	57.8	40.9	59.3	54.1	53.3
R-CNN (EB)	68.3	63.8	46.5	29.6	27.8	56.6	54.5	63.3	26.2	45.3	37.8	66.4	57.4	64.3	50.5	24.7	52.8	39.3	50.4	52.8	48.9
Ours	<b>76.8</b>	<b>71.2</b>	<b>61.2</b>	<b>45.1</b>	<b>35.9</b>	<b>62.5</b>	<b>60.9</b>	<b>75.5</b>	<b>31.3</b>	<b>58.3</b>	39.4	<b>73.8</b>	<b>68.6</b>	<b>73.1</b>	<b>60.9</b>	<b>33.1</b>	<b>59.1</b>	<b>41.0</b>	<b>66.3</b>	<b>57.6</b>	<b>57.7</b>



Sample detection results on the VOC 2012 dataset.