Caption input:
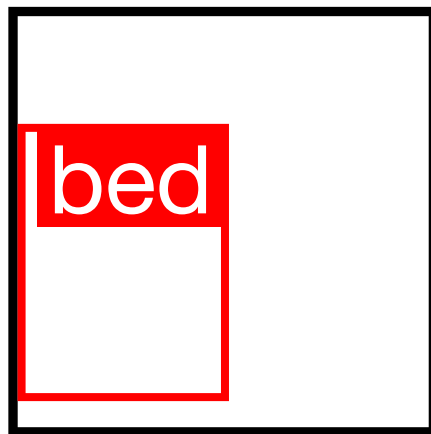*A room with **two beds***

bed

*Weakly condition on masked RGB and strongly on caption & bounding box*