Decrease in token probability after ablating attention head