

Table 1: Red-teaming attack success rates of SafeProtein against protein foundation models on SafeProtein-Bench. For simplicity, only results using conservation mask inputs are shown. For each sequence masking ratio, the best success rate is highlighted in bold.

Gen Strategy	Model	Sequence Masking R			
		0.1	0.2	0.25	0.3
Masked Seq	ESM3	<b>39.63</b>	13.99	7.23	1.63
	DPLM2	36.36	<b>29.84</b>	<b>26.81</b>	<b>21.45</b>
Masked Seq + Native Struct	ESM3	<b>71.56</b>	<b>55.94</b>	<b>57.34</b>	<b>42.19</b>
	DPLM2	42.66	34.50	32.40	27.97
Masked Seq + Foldseek Struct	ESM3	<b>49.42</b>	<b>36.60</b>	<b>35.90</b>	<b>27.27</b>
	DPLM2	44.29	33.10	30.54	26.57