

Table 1: Red-teaming attack success rates of SafeProtein against ESM3 on SafeProtein-Bench. Results are reported for all masking strategies. For each sequence masking ratio, the best success rate is highlighted in bold.

Gen Strategy	Masking Strategy	Sequence Masking Rate					
		0.1	0.2	0.25	0.3	0.4	0.5
Strategy1	Conservation	<b>39.63</b>	<b>13.99</b>	<b>7.23</b>	1.63	0.00	0.00
	Random	19.35	6.53	6.53	<b>4.20</b>	<b>3.33</b>	0.00
	Tail	5.83	0.93	1.17	0.70	0.00	0.00
Strategy2	Conservation	<b>71.56</b>	<b>55.94</b>	<b>57.34</b>	<b>42.19</b>	<b>39.00</b>	<b>35.00</b>
	Random	44.29	12.35	14.69	8.63	8.00	8.00
	Tail	34.03	7.69	9.79	3.73	4.00	4.00
Strategy3	Conservation	<b>49.42</b>	<b>36.60</b>	<b>35.90</b>	<b>27.27</b>	<b>22.00</b>	<b>21.00</b>
	Random	18.88	5.36	7.93	4.20	3.33	3.33
	Tail	6.99	1.87	2.10	1.17	2.00	2.00
Strategy4	Conservation	<b>72.49</b>	<b>63.64</b>	<b>64.10</b>	<b>46.85</b>	<b>43.00</b>	<b>42.00</b>
	Random	52.68	18.65	22.14	11.66	11.00	11.00
	Tail	42.89	10.72	15.15	7.93	10.00	10.00
Strategy5	Conservation	<b>75.06</b>	<b>75.06</b>	<b>74.36</b>	<b>72.26</b>	<b>72.00</b>	<b>72.00</b>
	Random	75.06	72.73	73.66	62.47	72.00	72.00
	Tail	74.13	65.73	66.43	52.21	66.00	66.00