- For instance: each number is treated as a separate token

- Replace them with a number token NUM
  - but: we want our language model to prefer

    $$p_{\text{LM}}(\text{I pay } 950.00 \text{ in May } 2007) > p_{\text{LM}}(\text{I pay } 2007 \text{ in May } 950.00)$$

  - not possible with number token

    $$p_{\text{LM}}(\text{I pay NUM in May NUM}) = p_{\text{LM}}(\text{I pay NUM in May NUM})$$

- Replace each digit (with unique symbol, e.g., @ or 5), retain some distinctions

  $$p_{\text{LM}}(\text{I pay } 555.55 \text{ in May } 5555) > p_{\text{LM}}(\text{I pay } 5555 \text{ in May } 555.55)$$