**a.** A model's explainability means how well its decision or output can be explained. This relies on explaining or understanding the model's reasoning behind its decisions, which factors affect its decision the most and how these decisions are derived from data. A model's explainability is crucial especially for critical applications such as regulatory systems or healthcare. Without understanding the models decisions, it would almost be impossible to debug underlying biases in data or ensure fairness.

**b.** Usually, especially in neural networks, more complex networks with more parameters and more hidden layers tend to perform better. However, the more complicated the model gets harder it becomes to understand the reasoning or dominant factors behind the model's decisions.

**c.** In the hypothetical scenario that we discussed in the hw2 FRQ's about the policing system of Minneapolis. Where a facial recognition system was applied to predict criminals. I think this is a great example where we would want to apply a more explainable model. As we can explain the model's decision more we can debug biases in the data with more certainty and help produce new data that has less bias than historical data.

Even though the model itself leads to higher inaccuracy as long as a supervisor can understand its decisions, the model's + supervisors decision can overall lead to higher accuracy. So, it could be generalized that when the decisions are highly critical, a more explainable model can be preferred.

**d.** I think LLM's and Chatbots such as ChatGPT or Gemini are great examples of where performance can be more beneficial than accuracy. One example application of these Chatbots are used to get summaries of complicated topics. A slightly wrong response to these questions can often be realized by the suspicious reader. Even in the worst case scenario when an LLM leads the reader to have a false understanding of the subject, knowledge can often be easily adjusted.