

Zero-Shot Model Performance Exploration Tests

```
graph TD; A[Zero-Shot Model Performance Exploration Tests] --> B[Model Forward]; A --> C[Dataset Forward]; B --> D[1. Answer Length Analysis<br/>▸ Are LMs capable of generating long answer spans?<br/>2. Sense Exploration<br/>▸ How good are LMs at detecting senses of key entity terms?<br/>3. Architecture Examination<br/>▸ Do variations on the same architecture (small v/s large v/s distilled, etc.) have an impact on performance?<br/>▸ Are bidirectional models better at this task than autoregressive models?<br/>4. Question category analysis<br/>▸ Quantitative breakdown of performance according to question type]; C --> E[1. (Dis)similarity between datasets<br/>▸ Quantitative measure using Force-Directed Algorithm<br/>2. Perplexity analysis<br/>▸ Is model performance correlated with dataset perplexity?<br/>3. Entity context coverage<br/>▸ How different are the contexts containing entities across datasets?];
```

Model Forward

1. Answer Length Analysis

- *Are LMs capable of generating long answer spans?*

2. Sense Exploration

- *How good are LMs at detecting senses of key entity terms?*

3. Architecture Examination

- *Do variations on the same architecture (small v/s large v/s distilled, etc.) have an impact on performance?*
- *Are bidirectional models better at this task than autoregressive models?*

4. Question category analysis

- *Quantitative breakdown of performance according to question type*

Dataset Forward

1. (Dis)similarity between datasets

- *Quantitative measure using Force-Directed Algorithm*

2. Perplexity analysis

- *Is model performance correlated with dataset perplexity?*

3. Entity context coverage

- *How different are the contexts containing entities across datasets?*