## 1029 A Taxonomy of Non-Strategic Microeconomics

1030

1050

1070

#### A.1 DECISIONS ON CONSUMPTION IN NON-STRATEGIC ENVIRONMENTS

We begin by characterizing the space of elements that test an agent's ability to optimally allocate 1031 their limited resources to goods and services they desire. In economics and decision theory, the 1032 most primitive approach to describing the preferences of decision-makers is to use a function that 1033 maps a set of possible choices to the agent's optimal choice within that set. Under a set of intuitive 1034 assumptions, such as transitivity (i.e., if bundle X is preferred to bundle Y, and Y is preferred to 1035 bundle Z, then X must be preferred to Z), it becomes possible to "rationalize" preferences by instead 1036 describing a utility function. This function assigns a real number to each bundle, and the agent selects 1037 the bundle with the highest utility. 1038

In this paper, we focus on these "rationalizable" preferences, where agent choice can be implemented 1039 as utility maximization constrained by prices and income. The solution to these consumer choice 1040 problems provides us with, among other things, individual demand functions, which describe the 1041 choice of each good or service as a function of prices and income. The individual demand functions 1042 for each good are essential when aggregating to the market demand in Consumer Goods Market 1043 Aggregation, which in turn is used to find the price in a non-strategic equilibrium. In addition, we test variations on the framework such as the agents ability to make tradeoffs between the quantity of 1045 goods they would need to be able to purchase for an increase in the amount of work they provide for 1046 a given wage (i.e., the elasticity of labor supply), as well as cases of choice under uncertainty where 1047 the agent is choosing between possible lotteries under rationalizability assumptions required for von 1048 Neuman expected utility. 1049

## **A.1.1** Properties of Utility Functions

In this section, we test the ability of the agent to use utility functions as a means to compare preferences 1051 over different "bundles" of goods or services. A key feature of economic reasoning in this context is 1052 for agents to consider how substitution between different goods in a bundle might achieve the same 1053 utility (i.e., map out the "indifference curves"). Key tests include correctly distinguishing between 1054 substitutes and complements in consumption, and calculating the marginal rate of substitution at a 1055 point on an indifference curve. This logic is essential for both agents acting as a planner as we will 1056 see in Appendix A.4 and when fulfilling the role of choice under budget and income constraints, in 1057 Deriving Demand 1058

Element A.1 (Marginal Utility). The ability to calculate marginal utility for different types of demand curves such as quasilinear, Cobb-Douglas, and Leontief.

Element A.2 (Diminishing Marginal Utility). The ability to recognize the role of diminishing marginal utility in consumption decisions and the role of achieving interior solutions.

Element A.3 (Marginal Rate of Substitution). The ability to calculate the marginal rate of substitution between two goods in a consumption decision.

Element A.4 (Tangency and the Marginal Rate of Substitution). The ability to calculate the marginal rate of substitution between two goods in a consumption decision at a given point in the budget constraint as tangent to the indifference curve.

1068 **Element A.5 (Substitutes and Complements).** The ability to distinguish between substitutes and complements in consumption decisions.

#### A.1.2 Deriving Demand

The module in this section tests an agent's ability to solve a constrained utility maximization problem

to derive a demand function—relying on the results of Properties of Utility Functions. We test the canonical classes of demand functions, check the duality of Marshallian demand and Hicksian

demand, and ask the agent to derive these demand functions from first principles.

Element A.6 (Derivation of Marshallian Demand). The ability to calculate the demand curve for a good given a utility function and a budget constraint.

Element A.7 (Derivation of Hicksian Demand from Expenditure Minimization). The ability to calculate the demand curve for a good given a utility function and a budget constraint.

Element A.8 (Duality of Hicksian Demand). The ability to recognize that Hicksian demand (expenditure minimization) is dual to maximization in Marshallian Demand.

#### 1081 A.1.3 Comparative Statics of Demand

- This module considers how agents reason about changes in prices or income, and their effects on the quantity of each good they would purchase. We test the classic law of demand, different types of goods (e.g., normal, inferior, and Giffen), and derive Engel curves from first principles. The key tests are to ensure the agent rationally responds to changes in relative prices, and investigate their substitution between goods in a bundle. In practice, these tests involve comparative statics of the argmax from the utility maximization of the previous section on Deriving Demand—i.e., using an Envelope theorem and perturbing prices or income.
- Element A.9 (Law of Demand). The ability to calculate the change in demand with the change in price for normal goods.
- Element A.10 (Price Elasticity of Demand). The ability to calculate the price elasticity of demand for a good given a utility function and a budget constraint.
- Element A.11 (Consumption Changes). The ability to change the relative expenditures on goods given changes in relative prices with ordinary or Giffen goods.
- 1095 **Element A.12 (Engel Curves).** The ability to calculate the Engel curve for a good given a utility 1096 function and a budget constraint.
- Element A.13 (Income Elasticity of Demand). The ability to calculate the income elasticity of demand for a good given a utility function and a budget constraint.

#### 1099 A.1.4 Labor Supply

- While the proceeding elements tested tradeoffs in choices of bundles with different goods, services (in Deriving Demand and over lotteries in Dynamic Consumption Decisions), often agents need to make a choice trading off between leisure and consumption. The elements in this module test an agent's ability to optimally make that tradeoff by balancing the consumption goods required to compensate for decreased leisure—which leads to the labor supply elasticity central to many branches of economics. Since goods must be purchased, agents will consider the relative wage from additional work compared to the price of goods. This leads us to be able to test an agent's ability to distinguish real from nominal prices.
- Element A.14 (Deriving Labor Supply). The ability to calculate the labor supply curve given specific preference parameterizations such as separable preferences or homothetic preferences.
- Element A.15 (Labor Supply Elasticity). The ability to calculate the elasticity of labor supply.
- Element A.16 (Marginal Rate of Substitution in Labor Supply). The ability to calculate the marginal rate of substitution between consumption and leisure in a labor supply decision.

#### 1113 A.1.5 Dynamic Consumption Decisions

- Individuals often face decisions about how to trade off more consumption today at the cost of additional debt and less consumption in the future, and how best to plan for consumption with various contingencies with the future is uncertain. Among other applications, this provides a formal model of how to best choose a mixture of financial assets—i.e., portfolios. Consequently, this subsection tests intertemporal consumption choices, optimal portfolio choice—which involves selecting a mix of assets that maximizes expected utility given the risks and returns associated with each asset. Understanding portfolio choice helps explain how consumers manage risk and make investment decisions, which is vital for financial planning and economic stability.
- Element A.17 (Price of Risk with Mean-Variance Utility). The ability to calculate the price of risk for a mean-variance utility function.
- Element A.18 (State-Contingent Consumption). The ability to calculate the optimal consumption given a utility function and a set of state-contingent consumption bundles.
- Element A.19 (Arbitrage). The ability to recognize and execute arbitrage opportunities given two goods and prices you can resell.

- Element A.20 (Optimal Portfolio Choice with Bid-Ask Spreads). The ability to calculate the optimal portfolio given bid-ask spreads.
- 1130 **Element A.21 (Exponential Discounting).** *The ability to exponentially discount future rewards or* 1131 *costs.*
- Element A.22 (Intertemporal Consumption Smoothing). The ability to calculate a smoothed consumption path and determine whether it is preferred to a non-smoothed path.

#### 1134 A.2 DECISIONS ON PRODUCTION IN NON-STRATEGIC ENVIRONMENTS

In the previous section, we derived how an agent facing a set of prices would choose the quantity 1135 demanded of each good or service to maximize their utility function. We also tested the amount of 1136 time that an agent might choose to wok (i.e., the quantity of labor supplied) given market wages-1137 where the agent trades off the additional goods they might purchase against the lost leisure time 1138 they must forgo. Here, we look at the other side of the market and test an agent's ability to operate 1139 1140 a production technology to maximize profits. Facing market prices for all production factors (e.g., wages and the capital) and the market price of the good or service they produce, the agent chooses the 1141 quantity of each factor of production and the total output. Parallel to DECISIONS ON CONSUMPTION 1142 IN NON-STRATEGIC ENVIRONMENTS, in Properties of Production Functions we first test general 1143 properties of production functions to ensure the agent can reason about substitution between factors, 1144 economies of scale in production, etc. Then in Deriving Factor Demand we solve the firms optimal 1145 profit maximization problem to determine the optimal choice of factors of production and output given a set of market prices. Finally, in Comparative Statics with Production we test the agents ability to reason about comparative statics on prices and their impact on factor demand and firm output. 1148

#### 1149 A.2.1 Properties of Production Functions

- Production functions in these environments take continuous inputs of each factor, which lets us test an agent's ability to conduct marginal thinking when choosing the composition of inputs. For example, by knowing the hourly wage of hiring an additional worker, the additional output the worker might produce using the particular production process, and the price they can sell the firm's output, they can decide whether hiring the additional worker is profitable. In the absence of prices, this section tests basic decision making of the agent for understanding substitution between factor of production, marginal products for each input, and the understanding of the returns to scale of a production process.
- Element A.23 (Marginal Products). The ability to calculate separate marginal products for a production function with multiple inputs (e.g., labor and capital).
- Element A.24 (Input Price Elasticity). The ability to calculate the responsiveness of output to a proportional change in a specific input's cost, holding all other inputs constant.
- Element A.25 (Output Elasticity). The ability to calculate the output elasticity of an input in a production function.
- Element A.26 (Elasticity of Substitution). *The ability to calculate the marginal elasticity of substitution between inputs in a production function.*
- Element A.27 (Diminishing Marginal Products). The ability to calculate the diminishing marginal products for a production function with multiple inputs.
- Element A.28 (Average and Marginal Costs). The ability to calculate average and marginal costs given a production function and input prices, and use them to determine scale.
- Element A.29 (Returns to Scale). The ability to determine the proportional change in output resulting from a proportional change in all inputs in a production function.

## A.2.2 Deriving Factor Demand

This module tests the agent's ability to act in the role of a profit maximizer in non-strategic situations where they take as given the price which they could sell goods they produce, and must pay for inputs to their production process at market rates (e.g., a competitive wage). Whereas in Deriving Demand the agent was solving a utility maximization problem subject to a budget constraint, here they solve a profit maximization problem constrained by a production function. We test decisions on the quantity and composition of inputs, and the quantity of output for canonical production functions such as Cobb-Douglas and Leontief production functions given the agent's understanding of production

- functions from Properties of Production Functions. The agent is asked to derive the factor demand
- functions from first principles from profit maximization and test their ability to reason with the dual
- cost-minimization formulation—analogous to the Hicksian vs. Marshallian demand of Deriving
- 1182 Demand

1191

- 1183 **Element A.30 (Profit Maximization).** The ability to calculate the optimal input bundle for a firm
- 1184 given a production function and input prices. Examples of given production functions: Cobb-Douglas,
- Leontief, Perfect Substitutes, CES production, CRS production, fixed costs.
- 1186 Element A.31 (Expenditure Minimization). The ability to calculate the optimal input bundle for a
- 1187 firm given a production function and input prices.
- 1188 Element A.32 (Duality of Profit Maximization and Expenditure Minimization). The ability to
- 1189 recognize that profit maximization is dual to expenditure minimization in production decisions and
- 1190 achieve consistent solutions.

#### A.2.3 Comparative Statics with Production

- This module considers how agents reason about changes in the prices at which they can sell their
- goods, as well as changes in the costs of producing those goods. In particular, we can test how
- this affects their optimal choice of inputs to their production process (e.g., how many people to
- hire or robots to lease). We test comparative statics on the prices of inputs to the production function, changes to the underlying production technology, and substitution between goods for classic
- function, changes to the underlying production technology, and substitution between goods for classic
- production functions such as Cobb-Douglas and Leontief. Analogous to the relationship between
- Deriving Demand and Comparative Statics of Demand, these tests involve comparative statics of the
- argmax from the profit maximization of Deriving Factor Demand—i.e., using an Envelope theorem
- and perturbing factor prices.
- 1201 Element A.33 (Price Elasticity of Supply). The ability to calculate the price elasticity of supply for
- a good given a production function and input prices.
- 1203 Element A.34 (Shephard's Lemma). The ability to calculate factor demands given a cost function
- using the derivatives with respect to prices.
- 1205 Element A.35 (Input Price Elasticity). The ability to calculate how the optimal input bundle changes
- with changes in input prices for a given production function.
- 1207 **Element A.36 (Total Factor Productivity).** The ability to calculate total factor productivity given a
- 1208 production function and input prices

### 1209 A.2.4 Dynamic Production Decisions

- While Deriving Factor Demand tested the ability of agents to make static (i.e., within-period) decisions
- on the mix of input factors to maximize profits, many producer problems are inherently dynamic.
- For example, we can test if an agent can optimally choose the amount of capital to purchase given
- 1213 forecasts of future consumer demand and prices or choose how much to adjust the labor force in
- cases when labor is difficult to relocate due to frictions such as hiring and firing costs. Finally, agents
- are tested on their ability to make optimal entry and exit decisions based on their forecasted profits in
- 1216 an evolving market.

1221

- 1217 Element A.37 (Dynamic Profit Maximization). The ability to calculate the optimal investment
- decision given a production function and input prices.
- 1219 Element A.38 (Entry and Exit Decisions). The ability to calculate the optimal entry and exit
- decisions given a production function and fixed costs.

#### A.3 DECISIONS IN MULTI-AGENT NON-STRATEGIC ENVIRONMENTS

- 1222 This setting tests the core logic of the relationship between supply-and-demand and prices, building
- on the tests of optimal behavior in appendix A.2 and appendix A.1. Economists refer to "general
- equilibrium" as the process where equilibrium prices and quantities emerge with a large number
- of non-strategic, price-taking market participants interact. Unlike the strategic models found in
- 1226 STEER, the assumption is that the market interactions that lead to this equilibrium occur through an
- unspecified process that clears markets (i.e., a "Walrasian auctioneer" or "invisible hand").
- In particular, for non-strategic settings, all market participants take prices as given and choose the
- quantity demanded or supplied in each market. For example, consumers jointly decide on the quantity

demanded of goods and services given relative prices, and the quantity of labor supplied given a wage.
Simultaneously, producers choose the quantity supplied of the good and the demand of each factor of production. With a large number of non-strategic market participants we can test the agents ability aggregate all of their supply and demand functions to calculate a market-level supply and demand. Finally, given the aggregated supply and demand functions for each market, we can test whether an agent can find the market clearing price where supply is equal to demand in equilibrium—given their

In this section, we organize by markets rather than by the role of a decision maker, as in the previous sections. For example, in the goods market we first ensure agents understand how individual demand functions from Deriving Demand aggregate to a market demand function for the good given a price, then that the agent understands how to aggregate the output from each producer at a given price from Comparative Statics with Production, and finally that the agent is able to calculate the price which would equate demand and supply and clear the market in a non-strategic setting. Factor markets are treated similarly.

Finally, given a system of equations that defines an equilibrium price we can perturb primitives (e.g., technological factors, distortions on decisions such as tax rates, or exogenous prices not determined in equilibrium) to see how the market clearing price would respond. That is an essential tool for agents to be able to reason about the impact of interventions and distortions in Appendix A.4.

#### A.3.1 Consumer Goods Market Aggregation

internal model of all the market participants.

1236

1248

1261

1277

The market clearing prices in general equilibrium arise from the separate market-level demand and supply curves, which sums the demand or supply across all market participants at a given price. Here we test the aggregation of demand functions derived from individual preferences, as in Deriving Demand and Comparative Statics of Demand, to a market demand function that summarizes the total quantity demanded across all agents at a given price. Central to the tests is to verify that the agent can aggregate the demands of market participants with heterogeneous preferences. On the other side of the market, we test if the agent can aggregate the "supply functions" resulting from the optimal choice of factors in Deriving Factor Demand and Comparative Statics with Production

Element A.39 (Aggregation of Consumer Demand). The ability to calculate the aggregate demand for a good given primitives of demand into expenditure shares.

Element A.40 (Aggregation of Offer Curve for the Good). The ability to calculate the aggregate supply of a good given primitives of supply into production functions.

#### A.3.2 Factor Market Aggregation

As with the case of the goods market in Consumer Goods Market Aggregation the market demand and supply for factors of production are essential to find the market clearing price. For example, we test whether the agent can aggregate the individual labor supply curve decisions from market participants who work at a particular wage, following Labor Supply into a market labor supply curve. On the other side of the market, we test whether the agent can aggregate the labor demand in Deriving Factor Demand from producers into a market labor demand curve. The same tests are essential for all factors of production, including capital.

Element A.41 (Aggregation of Labor Demand). The ability to calculate the aggregate demand for labor given primitives of demand into expenditure shares.

Element A.42 (Aggregation of Capital Demand). The ability to calculate the aggregate demand for capital given primitives of demand into expenditure shares.

Element A.43 (Aggregation of Labor Supply). The ability to calculate the aggregate supply of labor given primitives of supply into production functions.

Element A.44 (Aggregation of Fixed Factor Supply). The ability to calculate the aggregate supply of capital given primitives of supply into production functions.

## A.3.3 Prices in Static Market Equilibrium

In this setting we test the agent's ability to reason about how prices emerge in non-strategic setting as a process of equating supply and demand, which in turn relies on their ability to aggregate those market demand functions from consumer and producer behavior.

More specifically, the core logic of general equilibrium is to find the equilibrium price by taking the aggregated demand and supply functions for each market and find the prices which would equate demand and supply. For example, the supply and demand functions for the good, as a function of the price, in Consumer Goods Market Aggregation, or the supply and demand functions for factors of production, as a function of factor prices in Factor Market Aggregation. This is done market by market, taking all other prices as given—which requires the agent reason through comparative statics of the solution to a system of equations while keeping everything else fixed.

Element A.45 (Find Equilibrium Price). The ability to calculate the equilibrium prices given a production function and a demand function.

Element A.46 (Factor Shares in Equilibrium). The ability to calculate the factor shares in a competitive equilibrium given a production function and input prices.

### A.3.4 Comparative Statics of Equilibrium Prices

1292

1303

1326

Here, we test whether agents can reason about how prices and allocations (e.g., labor, capital, and goods) would respond to changes in the environment. The canonical tests are to see how changes in model primitives (e.g., productivity of the production process) or exogenous forces from outside the model (e.g., impact of weather), change the equilibrium price and allocations of labor, capital, etc. that would clear the market and equate demand and supply.

Element A.47 (Comparative Statics with Total Factor Production Shocks). The ability to calculate how equilibrium prices change with changes in input prices for a Cobb-Douglas production function.

Element A.48 (Comparative Statics with Inelastic or Perfectly Elastic Supply). The ability to calculate how equilibrium prices change with changes in input prices for a production function with inelastic or perfectly elastic supply.

#### A.4 EVALUATING EQUILIBRIA AND EXTERNALITIES

In Decisions in Multi-Agent Non-Strategic Environments, we tested an agents ability to reason about equilibrium prices and quantities arising from supply and demand decisions in a non-strategic setting. Although preferences were reflected in the underlying supply and demand functions themselves (i.e., utility maximization in the consumption decisions of Decisions on Consumption in Non-Strategic Environments and profit maximization in the production decisions of Decisions on Production in Non-Strategic Environments), the equilibria in Decisions in Multi-Agent Non-Strategic Environments do not necessarily reflect broader social preferences.

However, we can still ask whether the resulting "allocations" (i.e., the physical goods produced and how they are distributed to individuals, the amount of hours worked, and the physical capital installed) from the "invisible hand" in Decisions in Multi-Agent Non-Strategic Environments compare to a alternative ways of allocating resources which may directly take social preferences into account. A central result of economics in non-strategic settings is that absent market imperfections and market power (i.e., when self-interested agents cannot directly manipulate prices because they are too small) the competitive equilibria of Decisions in Multi-Agent Non-Strategic Environments typically yields the same allocations a benevolent planner might choose.

In this section, we consider how a social planner would evaluate the underlying welfare, efficiency, and inequality that comes about in non-strategic equilbria with prices derived from equating supply and demand. This leads to testing the ability of the agent to evaluate Pareto efficiency, consider the welfare theorems, evaluate Pigouvian externalities, and weigh the welfare impact of various market interventions which change the equilibria derived in DECISIONS IN MULTI-AGENT NON
STRATEGIC ENVIRONMENTS.

#### A.4.1 Welfare and Decentralization

In this section, we test whether the agent can determine cases where the the competitive equilibrium they calculate would yield the same distribution of resources and consumer welfare as that of a benevolent social planner directly making the consumption and production decisions of all agents directly (also known as the "Welfare Theorems"). In cases where the supply-and-demand relationships lead to the same results as those of a planner, the competitive equilibrium and its prices are said to

- "decentralize" the problem of a social planner. We then test that the agent recognizes cases where the
- welfare theorems fail, and can calculate the degree of welfare loss due to the distortions.
- Element A.49 (First Welfare Theorem). The ability to recognize that a competitive equilibrium is
- 1335 Pareto efficient.
- 1336 Element A.50 (Second Welfare Theorem). The ability to recognize that any Pareto efficient alloca-
- tion can be achieved as a competitive equilibrium with prices.
- 1338 Element A.51 (Consumer Surplus). The ability to calculate consumer surplus given a demand
- 1339 curve and a price.
- 1340 Element A.52 (Producer Surplus). The ability to calculate producer surplus given a supply curve
- 1341 and a price.
- 1342 Element A.53 (Efficient Surplus). The ability to calculate the total surplus in a competitive
- 1343 equilibrium and recognize that it is maximized in the competitive equilibrium.
- Element A.54 (Deadweight Loss of a Monopoly). The ability to calculate the deadweight loss of a
- monopoly given a demand curve and a supply curve.

#### 1346 A.4.2 Welfare Analysis of Market Equilibrium

- 1347 In this section, we focus on the agent's ability to evaluate welfare implications of various forms of
- market equilibrium, particularly how different policies and distortions impact overall efficiency and
- resource allocation. The agent is tested on their understanding of how different interventions—such
- as taxes, subsidies, and price controls—affect welfare outcomes, and their ability to distinguish
- between distortionary and non-distortionary policies.
- 1352 **Element A.55 (Identify Non-Distortionary Taxes).** The ability to identify taxes which do not distort
- 1353 the allocation of resources.
- 1354 Element A.56 (Irrelevance of Tax Incidence). The ability to recognize that the incidence of a tax
- does not depend on who is legally responsible for paying the tax.
- 1356 Element A.57 (Labor Supply Distortions). The ability to determine the extent that labor taxes will
- 1357 distort labor supply and change aggregates and prices.
- 1358 Element A.58 (Capital Market Distortions). The ability to identify that taxing a fixed factor is
- 1359 non-distortionary, but distorts with dynamic accumulation.

## 1360 B Mitigating Data Contamination with auto-STEER

- Data contamination, where training data inadvertently includes information from test sets, poses
- significant challenges in machine learning, leading to overestimated model performance and compro-
- mised generalization capabilities. To address this, we implemented a structured dataset generation
- methodology incorporating human oversight, controlled data generation, and style transfer techniques.
- This appendix details our approach and its alignment with best practices in the literature.
- The auto-STEER methodology provides a systematic approach to generating datasets that mitigates
- the risk of data contamination, ensuring the integrity of benchmarks and the validity of results. Below,
- we outline the key aspects of auto-STEER that address this issue:

#### **B.1** Challenging Models with Rephrasings:

- Rephrasings are known to cause significant variance in model performance, as demonstrated in the
- 1371 GSM-Symbolic dataset [32] and other studies [e.g., 61] [49] highlighting how syntactic or stylistic
- changes can challenge generalization. In Appendix G, we also show that much of the observed
- variance in LLM performance arises from these rephrasings, underscoring their role in robust
- evaluations. auto-STEER leverages this phenomenon to craft diverse rephrased questions that test
- beyond rote learning.

1369

### 1376 B.2 Dynamic Question Generation:

- 1377 auto-STEER generates new questions through a structured process that balances diversity and con-
- sistency. Questions are systematically rephrased or style-transferred to ensure they are different

- enough from the original templates to prevent memorization while retaining the same core meaning. 1379
- This approach reduces the risk of overlap with pre-trained data while preserving the focus of the 1380
- assessment. 1381
- The rapid advancement of large language models necessitates benchmarks that can evolve just as 1382
- quickly. To address this, auto-STEER incorporates a user interface that allows users to regenerate 1383
- entire datasets with minimal effort. By modifying domains, seeds, or even resampling numerical 1384
- values, users can quickly produce an entirely new dataset with minimal effort. This adaptability 1385
- ensures that benchmarks remain fresh and resistant to contamination as models advance. 1386

#### **Technical Descriptions of Metrics** $\mathbf{C}$ 1387

#### 1388 C.1 Accuracy.

- Accuracy is the most broadly used metric for evaluating LLMs. We define accuracy metrics as metrics 1389
- that only look at the top token that the LLM outputs. 1390

#### C.1.1 Exact-match accuracy 1391

- This is the fraction of questions answered correctly. In the FTQA format, we deem a LLM's response 1392
- correct if its final answer, when rounded to the same number of significant figures as the model's 1393
- output, matches the correct answer. This ensures that models are not rewarded for being vaguely 1394
- correct at low precision, while also penalizing overprecision: if an LLM reports more significant 1395
- figures than necessary and is incorrect, that discrepancy is treated as an error. This evaluation aligns 1396
- correctness with both the accuracy and confidence implied by the LLM's numerical output. 1397

#### C.1.2 Normalized accuracy 1398

- Elements can differ in their number of multiple choice options, leading to differences in the exact-1399
- match accuracy of random guessing. We can compensate for this by reporting the gap between the 1400
- 1401
- LLM's exact-match accuracy and random guessing  $\boxed{\textbf{G}}$ . We compute normalized accuracy for an element as follows:  $\sum_{i=1}^{N} a_i(t_i) \frac{1-a_i(t_i)}{|O_i|-1}$ , where  $t_i$  is the top token the LLM outputs for question 1402
- $i, a_i$  is the indicator describing whether the top token is correct or not, N the number of questions in 1403
- the element, and  $|O_i|$  the number of options in the question. In other words, normalized accuracy 1404
- rewards an LLM with 1 point for every correct answer and penalizes an LLM by 1 over the number 1405
- 1406 of options minus 1 for each incorrect answer.

#### **C.2** Calibration 1407

- It can also be useful to understand how confident an LLM is in its responses and the extent to which 1408
- these confidence levels align with accuracy.

#### C.2.1 Expected calibration error 1410

- We follow Liang et al. [28] and Raman et al. [38] in measuring the confidence of an LLM's response 1411
- and computing the expected calibration error [ECE; 34 17]. ECE measures how closely the probabil-
- ity an LLM assigns to its top answer matches the actual probability of the correct answer, which in
- our case is 1. ECE first splits the data into M equally spaced bins, where each bin contains the proba-1414
- bilities the model assigned to their top token in that range: e.g., let  $p^{\max}$  be the set of most probable tokens for each question then if M=2, then the first bin  $B_1=\{p\mid p\in p^{\max} \text{ and } p\in [0,0.5]\}$ . It is then defined as  $\sum_{i\in[M]}|B_i|/N\cdot|\operatorname{acc}(B_i)-\operatorname{conf}(B_i)|$ , where  $\operatorname{conf}(B_i)$  is the average probability 1415
- 1416
- 1417
- the LLM assigned to its top token in bin  $B_i$ , and  $acc(B_i)$  denotes the exact-match accuracy in bin 1418
- $B_i$ . We allow users to choose the number of bins, however, we set M=10 uniformly spaced over 1419
- the interval [0, 1] as is standard. 1420

#### 421 C.2.2 Brier score

The Brier Score of an element is defined as

$$\sum_{i=1}^{N} \frac{1}{|O_i|} \sum_{o \in O_i} (p_i(o) - a_i(o))^2,$$

where  $p_i(o)$  is the probability the LLM assigns to option o in question i. Thus, if an LLM is overly confident in an incorrect answer (e.g., assigns a probability of 0.9 to a wrong option), the Brier Score will penalize it more heavily.

#### 1425 C.2.3 Expected probability assignment

EPA measures how much probability mass an LLM assigns to the correct answer option out of all possible options. It is defined as:  $1/N \sum_{i \in [N]} p_i^*$ , where  $p_i^*$  is the predicted probability that the LLM assigns to the correct option for question i.

## 1429 D Technical Descriptions of Functional Families

In this section, we describe the functional forms that we use in testing economic concepts. Each can be applied to *consumer* problems (as utility functions) or *producer* problems (as production functions). We highlight the canonical mathematical form and note any technical differences in interpretation when modeling consumers versus producers.

#### 1434 D.1 Cobb-Douglas

The *Cobb-Douglas* functional form is one of the most frequently used due to its tractable properties and partial elasticities interpretation. Suppose there are n goods (or inputs). For a producer with input vector  $\mathbf{x} = (x_1, x_2, \dots, x_n)$ , a typical Cobb-Douglas production function can be written as:

$$f(\mathbf{x}) = A x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_n^{\alpha_n},$$

where A>0 is a scale parameter and each  $\alpha_i\geq 0$ . For a consumer's utility function, the same functional family looks like:

$$u(\mathbf{q}) = q_1^{\beta_1} q_2^{\beta_2} \cdots q_n^{\beta_n},$$

where  $\mathbf{q} = (q_1, q_2, \dots, q_n)$  are quantities of goods consumed, and  $\beta_i \geq 0$ . Economically,  $\alpha_i$  (or  $\beta_i$ ) often reflect the relative importance (or expenditure share) of each input (or good).

#### 1442 D.2 Leontief

A *Leontief* functional form encodes strict complementarity. A producer's Leontief production function is described as:

$$f(\mathbf{x}) = \min \left\{ \frac{x_1}{a_1}, \frac{x_2}{a_2}, \dots, \frac{x_n}{a_n} \right\},\,$$

where each  $a_i > 0$  captures a fixed proportion in which inputs must be combined. For a consumer, their Leontief utility function is of the form:

$$u(\mathbf{q}) = \min \left\{ \frac{q_1}{\gamma_1}, \frac{q_2}{\gamma_2}, \dots, \frac{q_n}{\gamma_n} \right\}.$$

This implies goods are perfect complements: the consumer gains utility only when goods are consumed in the specific ratio  $\gamma_1:\gamma_2:\ldots:\gamma_n$ . In production, perfect complementarity imposes that a shortage of any one input strictly limits total output.

#### D.3 Linear

1450

The *linear* family is the simplest and assumes perfect substitutability. For a producer, the linear production function with inputs **x** takes the form:

$$f(\mathbf{x}) = b_1 x_1 + b_2 x_2 + \dots + b_n x_n,$$

where  $b_i \ge 0$ . This means each input contributes additively (and independently) to output. A consumer's linear utility function with goods q is:

$$u(\mathbf{q}) = \theta_1 q_1 + \theta_2 q_2 + \dots + \theta_n q_n$$

where  $\theta_i > 0$  captures the marginal utility for good i. In both contexts, linear forms imply a constant rate of technical (or preferential) substitution, reflecting strong substitutability among inputs (or goods).

#### 458 D.4 Non-Linear

Beyond the classic forms, we use low-degree polynomial specifications for both production and utility functions. Such functions can capture nonlinearities without resorting to strictly Cobb-Douglas or other functional families. For a producer with inputs x, a second-degree polynomial looks like:

$$f(\mathbf{x}) = \alpha_0 + \sum_{i=1}^n \alpha_i x_i + \sum_{i=1}^n \sum_{j=1}^n \alpha_{ij} x_i x_j,$$

where  $\alpha_0, \alpha_i, \alpha_{ij}$  are specified parameters. Similarly, for a consumer's utility function, we write:

$$u(\mathbf{q}) = \beta_0 + \sum_{i=1}^n \beta_i q_i + \sum_{i=1}^n \sum_{j=1}^n \beta_{ij} q_i q_j.$$

These forms allow more flexible curvature than linear or Cobb-Douglas, but remain computationally tractable for low degrees. The main difference in interpretation is that producers typically focus on return-to-scale questions or marginal products, while consumers focus on marginal utilities and marginal rates of substitution.

#### 7 E Examples

#### E.1 GPT-40 Aggregation of Consumer Demand (Shown)

#### **User Message:**

Q: As the owner of a premium bedding brand, I cater to a diverse group of customers. I've discovered that there are 1204 potential buyers for my luxury bedding products. Out of these, 820 shoppers in urban markets each follow the demand function 18.1-1.51P, while 384 shoppers in suburban markets each adhere to the demand function 75.44-8.68P. How can I calculate the total demand for my bedding at a price point of 8.4?

- A. 2779.36
- B. 1832.01
- C. 5411.87 (Correct Answer)
- D. 5050.03

Let's think step by step. Explain your reasoning.

#### **GPT-40 Message:**

To calculate the total demand for your luxury bedding products at a price point of 8.4, we need to evaluate the demand functions for both urban and suburban markets separately and then sum the results.

### Step 1: Calculate the demand in the urban market

The demand function for the urban market is given by:

$$Q_u = 18.1 - 1.51P$$

Substitute P = 8.4 into the urban demand function:

$$Q_u = 18.1 - 1.51 \times 8.4$$

$$Q_u = 18.1 - 12.684$$

$$Q_u = 5.416$$

Since there are 820 shoppers in the urban market, the total demand in the urban market is:

Total Urban Demand = 
$$820 \times 5.416$$

Total Urban Demand = 4441.12

#### Step 2: Calculate the demand in the suburban market

The demand function for the suburban market is given by:

$$Q_s = 75.44 - 8.68P$$

Substitute P = 8.4 into the suburban demand function:

$$Q_s = 75.44 - 8.68 \times 8.4$$

$$Q_s = 75.44 - 72.912$$

$$Q_s = 2.528$$

Since there are 384 shoppers in the suburban market, the total demand in the suburban market is:

Total Suburban Demand =  $384 \times 2.528$ 

Total Suburban Demand = 970.88

#### Step 3: Calculate the total demand

Now, sum the total demand from both markets:

Total Demand = Total Urban Demand + Total Suburban Demand

Total Demand = 
$$4441.12 + 970.88$$

Total Demand 
$$= 5412$$

Therefore, the total demand for your bedding at a price point of 8.4 is approximately 5411.87.

The correct answer is:

C. 5411.87

1470

#### 471 E.2 GPT-40 Aggregation of Consumer Demand (Hidden)

#### **User Message:**

Q: As the owner of a premium bedding brand, I cater to a diverse group of customers. I've discovered that there are 1204 potential buyers for my luxury bedding products. Out of these, 820 shoppers in urban markets each follow the demand function 18.1-1.51P, while 384 shoppers in suburban markets each adhere to the demand function 75.44-8.68P. How can I calculate the total demand for my bedding at a price point of 8.4?

Let's think step by step. Explain your reasoning.

## **GPT-40 Message:**

To calculate the total demand for your luxury bedding at a price point of 8.4, we need to evaluate the demand functions for both urban and suburban markets separately and then sum the results. Here's a step-by-step breakdown:

- 1. Identify the demand functions:
- Urban market demand function:  $Q_u = 18.1 1.51P$
- Suburban market demand function:  $Q_s = 75.44 8.68P$
- 2. Substitute the price point (P = 8.4) into each demand function:
- For the urban market:

$$Q_u = 18.1 - 1.51 \times 8.4$$

Calculate the value:

$$Q_u = 18.1 - 1.51 \times 8.4 = 18.1 - 12.684 = 5.416$$

- For the suburban market:

$$Q_s = 75.44 - 8.68 \times 8.4$$

Calculate the value:

$$Q_s = 75.44 - 8.68 \times 8.4 = 75.44 - 72.912 = 2.528$$

3. Sum the demands from both markets to get the total demand:

$$Q_{total} = Q_u + Q_s = 5.416 + 2.528 = 7.944$$

Therefore, at a price point of 8.4, the total demand for your luxury bedding products from both urban and suburban markets is approximately 7.944 units.

1473

#### **User Message:**

- A. 2779.36
- B. 1832.01
- C. 5411.87 (Correct Answer)
- D. 5050.03

Answer by writing the option letter corresponding to the correct option. WRITE ONLY A SINGLE LETTER.

A:

1474

#### **GPT-40 Message:**

В

1475

#### 1476 E.3 GPT-40 Profit Maximization

#### **User Message:**

Q: John manages a company that produces premium leather wallets. The production output depends on the capital (K) he invests, represented by the function  $3.97K^{0.45}$ . At the moment, the wallets are sold at a price of 3.64, and John has an initial capital of  $K_1=3.3$ . He needs to decide how much additional capital to invest for the next day's market. John is aware of the potential price distribution for his wallets for the next day, represented by price 2.7 with probability 0.24, price 2.14 with probability 0.42, price 3.52 with probability 0.34. There is also a cost associated with increasing his capital, given by  $(K_1-K_2)^2$ , and he applies a discount factor of 0.37 to future revenues. To maximize his profits, how much additional capital should John invest?

- A. 0.44
- B. 0.14
- C. 1.1
- D. 0.09

Let's think step by step. Explain your reasoning.

#### **GPT-40 Message:**

To determine how much additional capital John should invest to maximize his profits, we need to consider the expected revenue from the different potential prices, the cost of increasing capital, and the discount factor applied to future revenues. Let's break this down step by step.

#### **Step 1: Calculate the Expected Revenue**

First, we need to calculate the expected price of the wallets for the next day using the given probabilities.

Expected Price = 
$$(2.7 \times 0.24) + (2.14 \times 0.42) + (3.52 \times 0.34)$$

Expected Price 
$$= 0.648 + 0.8988 + 1.1968 = 2.7436$$

Step 2: Determine the Production Function The production function is given by  $3.97K^{0.45}$ . We need to calculate the production output for different levels of capital K.

### Step 3: Calculate the Revenue for Different Capital Levels

Let's denote the initial capital as  $K_1 = 3.3$ . We need to consider the additional capital  $\Delta K$  that John might invest. The new capital  $K_2$  will be  $K_1 + \Delta K$ .

The revenue function for a given capital K is:

Revenue = Expected Price  $\times$  Production Output

Revenue = 
$$2.7436 \times 3.97K^{0.45}$$

#### **Step 4: Calculate the Cost of Increasing Capital**

The cost of increasing capital is given by  $(K_1 - K_2)^2$ . Since  $K_2 = K_1 + \Delta K$ , the cost becomes:

$$Cost = (\Delta K)^2$$

#### **Step 5: Calculate the Discounted Future Revenue**

John applies a discount factor of 0.37 to future revenues. Therefore, the discounted future revenue is:

Discounted Revenue = 
$$0.37 \times \text{Revenue}$$

## **Step 6: Calculate the Net Profit**

The net profit is the discounted revenue minus the cost of increasing capital:

Net Profit = 
$$0.37 \times (2.7436 \times 3.97 K_2^{0.45}) - (\Delta K)^2$$

#### **GPT-40 Message (Continued):**

#### Step 7: Evaluate the Net Profit for Different Values of $\Delta K$

We need to evaluate the net profit for the given options of additional capital investment:

Option A:  $\Delta K = 0.44$ 

$$K_2 = 3.3 + 0.44 = 3.74$$

Production Output =  $3.97 \times 3.74^{0.45}$ 

Revenue =  $2.7436 \times 3.97 \times 3.74^{0.45}$ 

Discounted Revenue =  $0.37 \times \text{Revenue}$ 

$$Cost = (0.44)^2 = 0.1936$$

Net Profit =  $0.37 \times \text{Revenue} - 0.1936$ 

Option B:  $\Delta K = 0.14$ 

$$K_2 = 3.3 + 0.14 = 3.44$$

Production Output =  $3.97 \times 3.44^{0.45}$ 

Revenue =  $2.7436 \times 3.97 \times 3.44^{0.45}$ 

Discounted Revenue =  $0.37 \times \text{Revenue}$ 

$$Cost = (0.14)^2 = 0.0196$$

Net Profit =  $0.37 \times \text{Revenue} - 0.0196$ 

Option C:  $\Delta K = 1.1$ 

$$K_2 = 3.3 + 1.1 = 4.4$$

Production Output =  $3.97 \times 4.4^{0.45}$ 

Revenue =  $2.7436 \times 3.97 \times 4.4^{0.45}$ 

Discounted Revenue =  $0.37 \times \text{Revenue}$ 

$$Cost = (1.1)^2 = 1.21$$

Net Profit =  $0.37 \times \text{Revenue} - 1.21$ 

Option D:  $\Delta K = 0.09$ 

$$K_2 = 3.3 + 0.09 = 3.39$$

Production Output =  $3.97 \times 3.39^{0.45}$ 

Revenue =  $2.7436 \times 3.97 \times 3.39^{0.45}$ 

Discounted Revenue =  $0.37 \times \text{Revenue}$ 

$$Cost = (0.09)^2 = 0.0081$$

Net Profit =  $0.37 \times \text{Revenue} - 0.0081$ 

#### **Step 8: Compare the Net Profits**

We need to compare the net profits for each option to determine which one maximizes John's profit. After calculating the net profits for each option, we find that:

- Option A:  $\Delta K = 0.44$  yields the highest net profit.

Therefore, the best option for John is to invest an additional capital of 0.44.

\*\*Answer: A. 0.44\*\*

1479

### **User Message:**

Answer by writing the option letter corresponding to the correct option. WRITE ONLY A SINGLE LETTER.

A:

1480

#### **GPT-40 Message:**

A

## 1482 F Figures

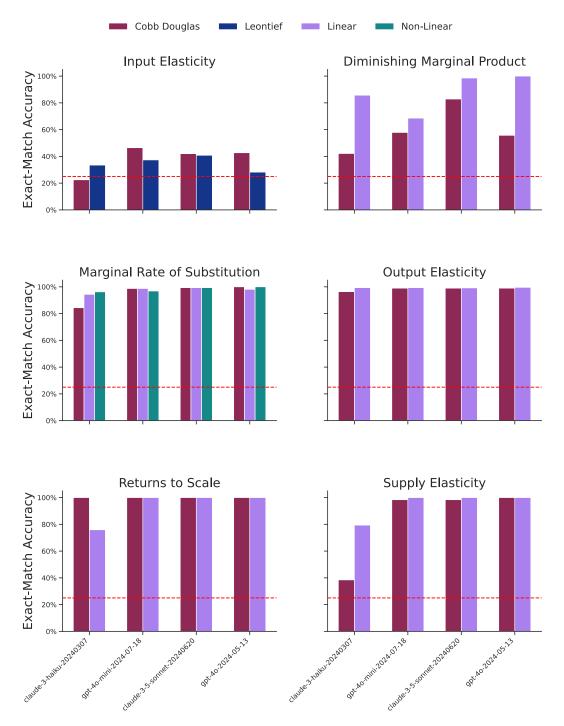


Figure 3: Exact-match accuracy of all closed-source models across six elements and four functional families (Cobb-Douglas, Leontief, Linear, and Non-Linear). The results demonstrate varying type robustness, with Cobb-Douglas being a generally more challenging functional family but not consistently harder for all elements. For instance, accuracy remains high for elements such as Output Elasticity and Marginal Rate of Substitution, even on the Cobb-Douglas functions, while elements like Input Price Elasticity and Returns to Scale show more variability across functional types. The red dashed line indicates the random guessing baseline for comparison.

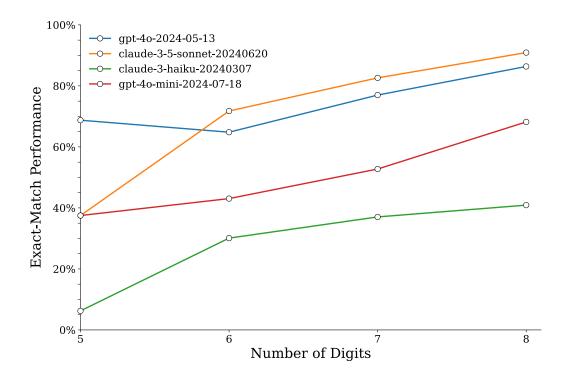


Figure 4: This figure depicts exact-match MCQA performance on the Aggregation of Consumer Demand element for the closed-source non-reasoning models against the number of digits of the correct answer.

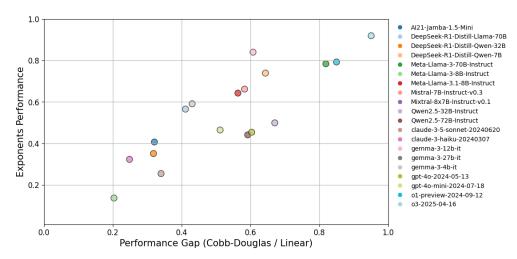


Figure 5: Scatter plot of calibrated performance on the Exponents element versus downstream performance gap across LLMs. The x-axis shows the gap calculated as the quotient between a LLM's accuracy on real-valued exponent-based (Cobb–Douglas) tasks and its accuracy on the linear version of those tasks for various downstream elements. The y-axis represents the LLM's performance on Exponents normalized by dividing by its average accuracy on the benchmark. Each point corresponds to a specific (LLM, downstream element) pair, with colors distinguishing different LLMs.

## **G** Analysis of Rephrasing Variance

semantically equivalent but syntactically varied prompts.

1494

To understand the role of question rephrasings in our dataset, we conducted an analysis of variance 1484 (ANOVA) on all other controllable features. These features include type, domain, and perspective. 1485 The goal of this analysis was to quantify the variance in LLM performance attributable to these 1486 features and, by exclusion, infer the contribution of rephrasings to the remaining unexplained variance. 1487 The results for the top-performing models, summarized in Table 2 through Table 6 indicate that 1488 the explained variance attributable to the controlled features is consistently low across all evaluated 1489 models. This leaves approximately 56% (for claude-3-5-sonnet) and up to 91% (for o1-preview) of 1490 the variance unexplained by the features included in the analysis. Given that question rephrasings 1491 are a systematic element of our dataset design and were not included as a feature in this analysis, 1492 we infer that the majority of this residual variance is due to differences in how models respond to 1493

Factor	Sum of Squares	Degrees of Freedom	F-Statistic	p-value
domain	7.5572	11.0000	3.0199	0.0823
perspective	3.4351	5.0000	3.0199	0.0823
CAR	21.1235	1.0000	92.8524	0.0000
element:type	178.6251	260.0000	3.0199	0.0823
Residual	2218.5356	9752.0000		
		R-squared		0.0941
		Adjusted R-squared		0.0907

Table 2: ANOVA Results for o1-preview-2024-09-12

Factor	Sum of Squares	Degrees of Freedom	F-Statistic	p-value
domain	0.2619	11.0000	0.0928	0.7607
perspective	0.1190	5.0000	0.0928	0.7607
0-CoT	229.7295	1.0000	895.0635	0.0000
CAR	141.3015	1.0000	550.5338	0.0000
element:type	47.4945	1995.0000	0.0928	0.7607
Residual	16015.5015	62399.0000		
		R-squared		0.3368
		Adjusted R-squared		0.3358

Table 3: ANOVA Results for gpt-4o-2024-05-13

Factor	Sum of Squares	Degrees of Freedom	F-Statistic	p-value
domain	0.4573	11.0000	0.1293	0.7191
perspective	0.2079	5.0000	0.1293	0.7191
0-Cot	34.0324	1.0000	105.8911	0.0000
CAR	159.1899	1.0000	495.3161	0.0000
element:type	82.9338	1995.0000	0.1293	0.7191
Residual	18642.5829	58006.0000		
		R-squared		0.2964
		Adjusted R-squared		0.2953

Table 4: ANOVA Results for gpt-4o-mini-2024-07-18

Factor	Sum of Squares	Degrees of Freedom	F-Statistic	p-value
domain	1.2389	11.0000	0.5004	0.6063
perspective	0.5032	5.0000	0.4472	0.5037
0-CoT	30.9486	1.0000	137.5118	0.0000
CAR	156.8135	1.0000	696.7577	0.0000
element:type	224.0990	1995.0000	0.4991	0.6071
Residual	18025.1941	80090.0000		
		R-squared		0.4436
		Adiusted R-sauared		0.4430

Adjusted R-squared
Table 5: ANOVA Results for claude-3-5-sonnet-20240620

Factor	Sum of Squares	Degrees of Freedom	F-Statistic	p-value
domain	1.4703	11.0000	0.3814	0.5369
perspective	0.6683	5.0000	0.3814	0.5369
0-Cot	0.2629	1.0000	0.7502	0.3864
CAR	0.4909	1.0000	1.4007	0.2366
element:type	266.6568	1995.0000	0.3814	0.5369
Residual	38842.9909	110826.0000		
		R-squared		0.2336
		Adjusted R-squared		0.2330

Adjusted R-squared
Table 6: ANOVA Results for claude-3-haiku-20240307

# 1495 H Models

Model Name	Model Card	Chat/ Instruction Tuned
Closed-Source		
OpenAI		
o3 o1-preview gpt-4o gpt-4o mini	_	√ √ √
Anthropic		
claude-3-5-sonnet claude-3-haiku <b>Open-Source</b>	-	<b>√</b> <b>√</b>
Google		
gemma-3-4b-it gemma-3-12b-it gemma-3-27b-it	google/gemma-3-4b-it google/gemma-3-12b-it google/gemma-3-27b-it	<b>√</b> <b>√</b> <b>√</b>
Qwen		
Qwen2.5-72B-Instruct Qwen2.5-7B-Instruct Qwen2.5-3B-Instruct	Qwen/Qwen2.5-72B-Instruct Qwen/Qwen2.5-7B-Instruct Qwen/Qwen2.5-3B-Instruct	<b>√</b> <b>√</b>
Qwen2.5-0.5B-Instruct Qwen2.5-math-7B-Instruct	Qwen/Qwen2.5-0.5B-Instruct Qwen/Qwen2.5-math-7B-Instruct	<b>,</b>
Qwen2.5-math-1.5B- Instruct	Qwen/Qwen2.5-math-1.5B-Instruct	<b>√</b> ✓
		Continued on next page

Commuea on next page

Model Card	Chat/Instruction Tuned
meta-llama/Meta-Llama-3-8B-Instruct	$\checkmark$
meta-llama/Meta-Llama-3-70B-Instruct	$\checkmark$
meta-llama/Meta-Llama-3.1-8B	×
meta-llama/Meta-Llama-3.1-70B-Instruct	$\checkmark$
mistralai/Mixtral-8x7B-Instruct-v0.1	$\checkmark$
mistralai/Mistral-7B-Instruct-v0.3	$\checkmark$
ai211abs/AI21-Jamba-1.5-Mini	×
	meta-llama/Meta-Llama-3-8B-Instruct meta-llama/Meta-Llama-3-70B-Instruct meta-llama/Meta-Llama-3.1-8B meta-llama/Meta-Llama-3.1-70B-Instruct  mistralai/Mixtral-8x7B-Instruct-v0.1 mistralai/Mistral-7B-Instruct-v0.3

Table 7: Overview of the open- and closed-source LLMs we evaluated. The table includes their names, their model card links, and whether they have been chat or instruction tuned. Models are grouped by family and sorted by parameter size, with non-chat-tuned models listed first within each group.

### I Extra Results

#### I.1 Performance on Elements Generated by Claude 3.5 Sonnet

To assess whether the performance on our dataset was influenced by the choice of the generation LLM, we re-generated three elements from scratch using claude-3-5-sonnet. We selected Find Equilibrium Price because it exhibited the largest performance gap between gpt-40 and claude-3-5-sonnet, Diminishing Marginal Products was chosen as a random element with slight performance variation across the models, and Price Elasticity of Demand served as a control where no significant differences were expected.

Figure 6 shows the exact-match performance of both models on these three elements. We found no significant differences in performance between any of the models.



Figure 6: Exact-match performance comparison between closed-source models on three that were generated by gpt-40 and claude-3-5-sonnet. The elements were selected based on the observed performance differences across models, with Price Elasticity of Demand serving as a control. The figure shows no significant differences in performance between the models on these elements. Note that the red dotted line signifies random guessing performance.

#### 1506 I.2 Intertemporal Consumption Smoothing

When optimizing intertemporal consumption, the consumer maximizes the discounted utility

$$\sum_{t=0}^{T} \beta^t u(c_t)$$

subject to the intertemporal budget constraint. The first-order condition for an optimum leads to the Euler equation:

$$u'(c_t) = \beta(1+r) u'(c_{t+1}).$$

For our purposes, we tested models using a constant relative risk aversion (CRRA) utility function.

1511 We used the following form:

$$u(c) = \frac{c^{1-\gamma}}{1-\gamma} \quad (\gamma \neq 1),$$

where  $\gamma$  is the coefficient of relative risk aversion. This function exhibits diminishing marginal utility, meaning that each additional unit of consumption adds less to utility when overall consumption is high compared to when it is low. Due to diminishing returns, an agent is motivated to smooth consumption over time, even with a discount factor of 1; spending too much in one period reduces the marginal utility in that period, while having too little in another period results in a steep loss of satisfaction.

1518 For CRRA utility, the Euler equation is given by:

$$u'(c_t) = \beta(1+r) u'(c_{t+1}),$$

1519 where

1527

1528

1530

1533

1534

1535

1536

1537

1538

1539

1540

1541

1542 1543

1544

1545

1546

$$u'(c) = c^{-\gamma}.$$

1520 Rearranging yields:

$$c_{t+1} = [\beta(1+r)]^{\frac{1}{\gamma}} c_t.$$

However, we found that LLMs would often use linear utility functions in their analysis. For a linear utility function, the Euler equation—which equates the marginal benefit of consuming today with that of consuming tomorrow—simplifies significantly. If u'(c) is constant (say, equal to 1), then aside from the effects of discounting and interest, there is no curvature-driven motive to adjust consumption levels across periods. The optimal allocation would then depend solely on the intertemporal budget constraint and the returns on savings.

#### I.3 Deadweight Loss

To conduct our error analysis, we ran all closed-source models on the free-text QA adaptation of the deadweight loss task. We began by inspecting a range of model outputs to identify distinct classes of errors that were common across responses. Once these error categories were established, we computed the answers corresponding to these errors and rescored the models based on whether their outputs were within 98% of either the correct answer or any of the answers derived from specific error assumptions. We also ensured that when there was any overlap in incorrect responses that we chose the closest one to the model's response. This approach allowed us to capture not only the frequency of correct outputs but also the systematic nature of the models' reasoning flaws. Below, we provide a detailed breakdown of the primary error types:

- Incorrect Base for Deadweight Loss Type 1: This error incorrectly substitutes  $P_e P_m$  (the difference between the competitive equilibrium price and the monopolist's price) in place of the correct term  $P_m MC(Q_m)$  (the difference between the monopolist's price and the marginal cost at the monopolist's quantity).
- Incorrect Base for Deadweight Loss Type 2: This error calculates the deadweight loss using
  the difference between the monopoly price and the competitive equilibrium price as the base
  of the triangle.
- Incorrect Base and Height Type 1: This error replaces the base of the DWL triangle  $(Q_e Q_m)$  with a miscalculated value for the equilibrium quantity and replaces the base with the Type 1 variant.

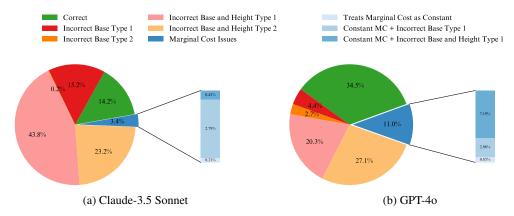


Figure 7: Error analyses of claude-3-5-sonnet and gpt-40 on the Deadweight Loss of a Monopoly element. In reds and oranges are failures due to incorrect computations of the deadweight loss area; in blue and further broken down are errors due to incorrectly interpreting the marginal cost. A more detailed description of what each error means can be found in Appendix I.3.

- Incorrect Base and Height Type 2: This error assumes that  $P_e = MC(Q_m)$ , leading to an incorrect height calculation where the DWL triangle's height becomes  $P_m P_e$ .
- Treating Marginal Cost as Constant: Instead of recognizing marginal cost as a function derived from the supply curve ( $MC(Q) = a_{\text{supply}} \times Q + b_{\text{supply}}$ ), several models treated marginal cost as a constant, often equal to the slope of the supply curve. This assumption led to errors in determining the monopolist's quantity and price, further propagating inaccuracies in the deadweight loss calculation.
- Combined Errors: A subset of models combined the two errors above, simultaneously using
  an incorrect formula for deadweight loss and assuming a constant marginal cost. This
  compounded error significantly reduced the likelihood of producing a correct answer and
  highlighted the systematic nature of the misunderstanding.

## 1558 J Section 3 Images



Figure 8: (Left) The hidden approach to 0-CoT: the model is given only the question and asked to explain its reasoning before being provided with options. (Right) The shown approach to 0-CoT: the model is presented with both the question and options before explaining its reasoning.

## 1559 K Web Application

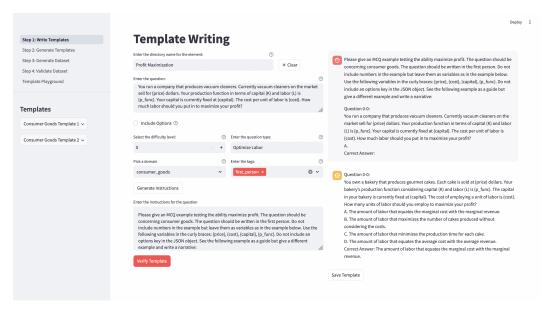


Figure 9: The web app user interface for template writing. This page includes fields for type, domain, grade level and tags (including perspectives). The right shows an example of template verification which uses a LLM to generate another template using the example seed.

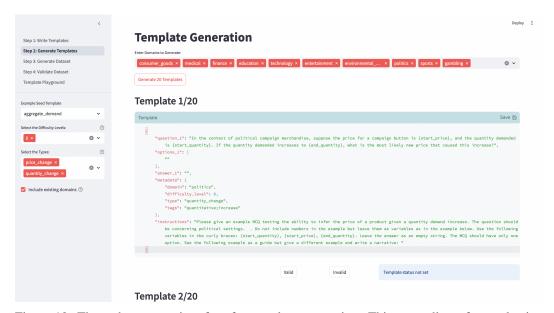


Figure 10: The web app user interface for template generation. This page allows for a selection of domains, and types for which templates will be generated using the available example seeds. Templates can then be verified and saved by the user.

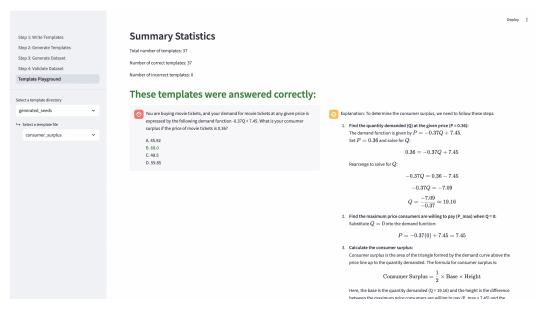


Figure 11: The web app user interface for template AI double-checking. This page instantiates and fills a set of question using a generated or example seed and then generates a response using an OpenAI model. The page also reports the number of questions answered correctly as well as the responses from the model.

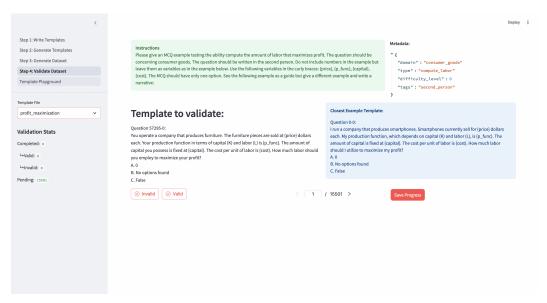


Figure 12: The web app user interface for template validation. This page displays all generated seeds returned by the model for manual validation.