

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

A TAXONOMY OF NON-STRATEGIC MICROECONOMICS

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 their limited resources to goods and services they desire. In economics and decision theory, the most primitive approach to describing the preferences of decision-makers is to use a function that maps a set of possible choices to the agent’s optimal choice within that set. Under a set of intuitive assumptions, such as *transitivity* (i.e., if bundle X is preferred to bundle Y , and Y is preferred to bundle Z , then X must be preferred to Z), it becomes possible to “rationalize” preferences by instead describing a utility function. This function assigns a real number to each bundle, and the agent selects the bundle with the highest utility.

In this paper, we focus on these “rationalizable” preferences, where agent choice can be implemented as utility maximization constrained by prices and income. The solution to these consumer choice problems provides us with, among other things, individual demand functions, which describe the choice of each good or service as a function of prices and income. The individual demand functions for each good are essential when aggregating to the market demand in **Consumer Goods Market 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 goods they would need to be able to purchase for an increase in the amount of work they provide for a given wage (i.e., the elasticity of labor supply), as well as cases of choice under uncertainty where the agent is choosing between possible lotteries under rationalizability assumptions required for von Neuman expected utility.

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 over different “bundles” of goods or services. A key feature of economic reasoning in this context is for agents to consider how substitution between different goods in a bundle might achieve the same utility (i.e., map out the “indifference curves”). Key tests include correctly distinguishing between substitutes and complements in consumption, and calculating the marginal rate of substitution at a point on an indifference curve. This logic is essential for both agents acting as a planner as we will see in **Appendix A.4** and when fulfilling the role of choice under budget and income constraints, in **Deriving Demand**

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.*

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.*

972 **Element A.8 (Duality of Hicksian Demand).** *The ability to recognize that Hicksian demand (expen-*
973 *diture minimization) is dual to maximization in Marshallian Demand.*

974

975

976

A.1.3 COMPARATIVE STATICS OF DEMAND

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

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.*

Element A.12 (Engel Curves). *The ability to calculate the Engel curve for a good given a utility 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.*

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

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.*

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

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.*

1026 **Element A.20 (Optimal Portfolio Choice with Bid-Ask Spreads).** *The ability to calculate the*
1027 *optimal portfolio given bid-ask spreads.*

1028 **Element A.21 (Exponential Discounting).** *The ability to exponentially discount future rewards or*
1029 *costs.*

1030 **Element A.22 (Intertemporal Consumption Smoothing).** *The ability to calculate a smoothed*
1031 *consumption path and determine whether it is preferred to a non-smoothed path.*
1032

1033 A.2 DECISIONS ON PRODUCTION IN NON-STRATEGIC ENVIRONMENTS 1034

1035 In the previous section, we derived how an agent facing a set of prices would choose the quantity
1036 demanded of each good or service to maximize their utility function. We also tested the amount of
1037 time that an agent might choose to work (i.e., the quantity of labor supplied) given market wages—
1038 where the agent trades off the additional goods they might purchase against the lost leisure time
1039 they must forgo. Here, we look at the other side of the market and test an agent's ability to operate
1040 a production technology to maximize profits. Facing market prices for all production factors (e.g.,
1041 wages and the capital) and the market price of the good or service they produce, the agent chooses the
1042 quantity of each factor of production and the total output. Parallel to DECISIONS ON CONSUMPTION
1043 IN NON-STRATEGIC ENVIRONMENTS, in **Properties of Production Functions** we first test general
1044 properties of production functions to ensure the agent can reason about substitution between factors,
1045 economies of scale in production, etc. Then in **Deriving Factor Demand** we solve the firm's optimal
1046 profit maximization problem to determine the optimal choice of factors of production and output
1047 given a set of market prices. Finally, in **Comparative Statics with Production** we test the agent's ability
1048 to reason about comparative statics on prices and their impact on factor demand and firm output.

1049 A.2.1 PROPERTIES OF PRODUCTION FUNCTIONS 1050

1051 Production functions in these environments take continuous inputs of each factor, which lets us test an
1052 agent's ability to conduct marginal thinking when choosing the composition of inputs. For example,
1053 by knowing the hourly wage of hiring an additional worker, the additional output the worker might
1054 produce using the particular production process, and the price they can sell the firm's output, they
1055 can decide whether hiring the additional worker is profitable. In the absence of prices, this section
1056 tests basic decision making of the agent for understanding substitution between factor of production,
1057 marginal products for each input, and the understanding of the returns to scale of a production process.

1058 **Element A.23 (Marginal Products).** *The ability to calculate separate marginal products for a*
1059 *production function with multiple inputs (e.g., labor and capital).*

1060 **Element A.24 (Input Price Elasticity).** *The ability to calculate the responsiveness of output to a*
1061 *proportional change in a specific input's cost, holding all other inputs constant.*

1062 **Element A.25 (Output Elasticity).** *The ability to calculate the output elasticity of an input in a*
1063 *production function.*

1064 **Element A.26 (Elasticity of Substitution).** *The ability to calculate the marginal elasticity of substi-*
1065 *tution between inputs in a production function.*

1066 **Element A.27 (Diminishing Marginal Products).** *The ability to calculate the diminishing marginal*
1067 *products for a production function with multiple inputs.*

1068 **Element A.28 (Average and Marginal Costs).** *The ability to calculate average and marginal costs*
1069 *given a production function and input prices, and use them to determine scale.*

1070 **Element A.29 (Returns to Scale).** *The ability to determine the proportional change in output*
1071 *resulting from a proportional change in all inputs in a production function.*

1072 A.2.2 DERIVING FACTOR DEMAND 1073

1074 This module tests the agent's ability to act in the role of a profit maximizer in non-strategic situations
1075 where they take as given the price which they could sell goods they produce, and must pay for inputs
1076 to their production process at market rates (e.g., a competitive wage). Whereas in **Deriving Demand**,
1077 the agent was solving a utility maximization problem subject to a budget constraint, here they solve a
1078 profit maximization problem constrained by a production function. We test decisions on the quantity
1079 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—analogueous to the Hicksian vs. Marshallian demand of [Deriving Demand](#).

Element A.30 (Profit Maximization). *The ability to calculate the optimal input bundle for a firm given a production function and input prices. Examples of given production functions: Cobb-Douglas, Leontief, Perfect Substitutes, CES production, CRS production, fixed costs.*

Element A.31 (Expenditure Minimization). *The ability to calculate the optimal input bundle for a firm given a production function and input prices.*

Element A.32 (Duality of Profit Maximization and Expenditure Minimization). *The ability to recognize that profit maximization is dual to expenditure minimization in production decisions and 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 production functions such as Cobb-Douglas and Leontief. Analogueous 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.

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.*

Element A.34 (Shephard’s Lemma). *The ability to calculate factor demands given a cost function using the derivatives with respect to prices.*

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.*

Element A.36 (Total Factor Productivity). *The ability to calculate total factor productivity given a production function and input prices*

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 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 an evolving market.

Element A.37 (Dynamic Profit Maximization). *The ability to calculate the optimal investment decision given a production function and input prices.*

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

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 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”).

1134 In particular, for non-strategic settings, all market participants take prices as given and choose the
1135 quantity demanded or supplied in each market. For example, consumers jointly decide on the quantity
1136 demanded of goods and services given relative prices, and the quantity of labor supplied given a wage.
1137 Simultaneously, producers choose the quantity supplied of the good and the demand of each factor of
1138 production. With a large number of non-strategic market participants we can test the agents ability
1139 aggregate all of their supply and demand functions to calculate a market-level supply and demand.
1140 Finally, given the aggregated supply and demand functions for each market, we can test whether an
1141 agent can find the market clearing price where supply is equal to demand in equilibrium—given their
1142 internal model of all the market participants.

1143 In this section, we organize by markets rather than by the role of a decision maker, as in the previous
1144 sections. For example, in the goods market we first ensure agents understand how individual demand
1145 functions from [Deriving Demand](#) aggregate to a market demand function for the good given a price,
1146 then that the agent understands how to aggregate the output from each producer at a given price from
1147 [Comparative Statics with Production](#), and finally that the agent is able to calculate the price which
1148 would equate demand and supply and clear the market in a non-strategic setting. Factor markets are
1149 treated similarly.

1150 Finally, given a system of equations that defines an equilibrium price we can perturb primitives (e.g.,
1151 technological factors, distortions on decisions such as tax rates, or exogenous prices not determined
1152 in equilibrium) to see how the market clearing price would respond. That is an essential tool for
1153 agents to be able to reason about the impact of interventions and distortions in [Appendix A.4](#)

1156 A.3.1 CONSUMER GOODS MARKET AGGREGATION

1158 The market clearing prices in general equilibrium arise from the separate market-level demand and
1159 supply curves, which sums the demand or supply across all market participants at a given price. Here
1160 we test the aggregation of demand functions derived from individual preferences, as in [Deriving](#)
1161 [Demand](#) and [Comparative Statics of Demand](#) to a market demand function that summarizes the total
1162 quantity demanded across all agents at a given price. Central to the tests is to verify that the agent
1163 can aggregate the demands of market participants with heterogeneous preferences. On the other side
1164 of the market, we test if the agent can aggregate the “supply functions” resulting from the optimal
1165 choice of factors in [Deriving Factor Demand](#) and [Comparative Statics with Production](#)

1166 **Element A.39 (Aggregation of Consumer Demand).** *The ability to calculate the aggregate demand*
1167 *for a good given primitives of demand into expenditure shares.*

1168 **Element A.40 (Aggregation of Offer Curve for the Good).** *The ability to calculate the aggregate*
1169 *supply of a good given primitives of supply into production functions.*

1172 A.3.2 FACTOR MARKET AGGREGATION

1174 As with the case of the goods market in [Consumer Goods Market Aggregation](#) the market demand
1175 and supply for factors of production are essential to find the market clearing price. For example,
1176 we test whether the agent can aggregate the individual labor supply curve decisions from market
1177 participants who work at a particular wage, following [Labor Supply](#) into a market labor supply curve.
1178 On the other side of the market, we test whether the agent can aggregate the labor demand in [Deriving](#)
1179 [Factor Demand](#) from producers into a market labor demand curve. The same tests are essential for all
1180 factors of production, including capital.

1181 **Element A.41 (Aggregation of Labor Demand).** *The ability to calculate the aggregate demand for*
1182 *labor given primitives of demand into expenditure shares.*

1183 **Element A.42 (Aggregation of Capital Demand).** *The ability to calculate the aggregate demand for*
1184 *capital given primitives of demand into expenditure shares.*

1185 **Element A.43 (Aggregation of Labor Supply).** *The ability to calculate the aggregate supply of*
1186 *labor given primitives of supply into production functions.*

1187 **Element A.44 (Aggregation of Fixed Factor Supply).** *The ability to calculate the aggregate supply*
of capital given primitives of supply into production functions.

1188 A.3.3 PRICES IN STATIC MARKET EQUILIBRIUM
1189

1190 In this setting we test the agent’s ability to reason about how prices emerge in non-strategic setting
1191 as a process of equating supply and demand, which in turn relies on their ability to aggregate those
1192 market demand functions from consumer and producer behavior.

1193 More specifically, the core logic of general equilibrium is to find the equilibrium price by taking the
1194 aggregated demand and supply functions for each market and find the prices which would equate
1195 demand and supply. For example, the supply and demand functions for the good, as a function of
1196 the price, in **Consumer Goods Market Aggregation**, or the supply and demand functions for factors
1197 of production, as a function of factor prices in **Factor Market Aggregation**. This is done market by
1198 market, taking all other prices as given—which requires the agent reason through comparative statics
1199 of the solution to a system of equations while keeping everything else fixed.

1200 **Element A.45 (Find Equilibrium Price).** *The ability to calculate the equilibrium prices given a*
1201 *production function and a demand function.*

1202 **Element A.46 (Factor Shares in Equilibrium).** *The ability to calculate the factor shares in a*
1203 *competitive equilibrium given a production function and input prices.*

1204

1205

1206 A.3.4 COMPARATIVE STATICS OF EQUILIBRIUM PRICES
1207

1208 Here, we test whether agents can reason about how prices and allocations (e.g., labor, capital, and
1209 goods) would respond to changes in the environment. The canonical tests are to see how changes in
1210 model primitives (e.g., productivity of the production process) or exogenous forces from outside the
1211 model (e.g., impact of weather), change the equilibrium price and allocations of labor, capital, etc.
1212 that would clear the market and equate demand and supply.

1213 **Element A.47 (Comparative Statics with Total Factor Production Shocks).** *The ability to calculate*
1214 *how equilibrium prices change with changes in input prices for a Cobb-Douglas production function.*

1215 **Element A.48 (Comparative Statics with Inelastic or Perfectly Elastic Supply).** *The ability to*
1216 *calculate how equilibrium prices change with changes in input prices for a production function with*
1217 *inelastic or perfectly elastic supply.*

1218

1219

1220 A.4 EVALUATING EQUILIBRIA AND EXTERNALITIES

1221

1222 In DECISIONS IN MULTI-AGENT NON-STRATEGIC ENVIRONMENTS, we tested an agents ability
1223 to reason about equilibrium prices and quantities arising from supply and demand decisions in a
1224 non-strategic setting. Although preferences were reflected in the underlying supply and demand
1225 functions themselves (i.e., utility maximization in the consumption decisions of DECISIONS ON
1226 CONSUMPTION IN NON-STRATEGIC ENVIRONMENTS and profit maximization in the production
1227 decisions of DECISIONS ON PRODUCTION IN NON-STRATEGIC ENVIRONMENTS), the equilibria
1228 in DECISIONS IN MULTI-AGENT NON-STRATEGIC ENVIRONMENTS do not necessarily reflect
broader social preferences.

1229

1230 However, we can still ask whether the resulting “allocations” (i.e., the physical goods produced and
1231 how they are distributed to individuals, the amount of hours worked, and the physical capital installed)
1232 from the “invisible hand” in DECISIONS IN MULTI-AGENT NON-STRATEGIC ENVIRONMENTS
1233 compare to a alternative ways of allocating resources which may directly take social preferences
1234 into account. A central result of economics in non-strategic settings is that absent market imperfec-
1235 tions and market power (i.e., when self-interested agents cannot directly manipulate prices because
1236 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.

1237

1238 In this section, we consider how a social planner would evaluate the underlying welfare, efficiency,
1239 and inequality that comes about in non-strategic equilibria with prices derived from equating supply
1240 and demand. This leads to testing the ability of the agent to evaluate Pareto efficiency, consider
1241 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.

1242 A.4.1 WELFARE AND DECENTRALIZATION

1243

1244 In this section, we test whether the agent can determine cases where the the competitive equilibrium
1245 they calculate would yield the same distribution of resources and consumer welfare as that of a
1246 benevolent social planner directly making the consumption and production decisions of all agents
1247 directly (also known as the “Welfare Theorems”). In cases where the supply-and-demand relationships
1248 lead to the same results as those of a planner, the competitive equilibrium and its prices are said to
1249 “decentralize” the problem of a social planner. We then test that the agent recognizes cases where the
1250 welfare theorems fail, and can calculate the degree of welfare loss due to the distortions.

1251 **Element A.49 (First Welfare Theorem).** *The ability to recognize that a competitive equilibrium is*
1252 *Pareto efficient.*

1253 **Element A.50 (Second Welfare Theorem).** *The ability to recognize that any Pareto efficient alloca-*
1254 *tion can be achieved as a competitive equilibrium with prices.*

1255 **Element A.51 (Consumer Surplus).** *The ability to calculate consumer surplus given a demand*
1256 *curve and a price.*

1257 **Element A.52 (Producer Surplus).** *The ability to calculate producer surplus given a supply curve*
1258 *and a price.*

1259 **Element A.53 (Efficient Surplus).** *The ability to calculate the total surplus in a competitive*
1260 *equilibrium and recognize that it is maximized in the competitive equilibrium.*

1261 **Element A.54 (Deadweight Loss of a Monopoly).** *The ability to calculate the deadweight loss of a*
1262 *monopoly given a demand curve and a supply curve.*

1263

1264 A.4.2 WELFARE ANALYSIS OF MARKET EQUILIBRIUM

1265

1266 In this section, we focus on the agent’s ability to evaluate welfare implications of various forms of
1267 market equilibrium, particularly how different policies and distortions impact overall efficiency and
1268 resource allocation. The agent is tested on their understanding of how different interventions—such
1269 as taxes, subsidies, and price controls—affect welfare outcomes, and their ability to distinguish
1270 between distortionary and non-distortionary policies.

1271 **Element A.55 (Identify Non-Distortionary Taxes).** *The ability to identify taxes which do not distort*
1272 *the allocation of resources.*

1273 **Element A.56 (Irrelevance of Tax Incidence).** *The ability to recognize that the incidence of a tax*
1274 *does not depend on who is legally responsible for paying the tax.*

1275 **Element A.57 (Labor Supply Distortions).** *The ability to determine the extent that labor taxes will*
1276 *distort labor supply and change aggregates and prices.*

1277 **Element A.58 (Capital Market Distortions).** *The ability to identify that taxing a fixed factor is*
1278 *non-distortionary, but distorts with dynamic accumulation.*

1279

1280

1281 B MITIGATING DATA CONTAMINATION WITH `AUTO-STEER`

1282

1283 Data contamination, where training data inadvertently includes information from test sets, poses
1284 significant challenges in machine learning, leading to overestimated model performance and compro-
1285 mised generalization capabilities. To address this, we implemented a structured dataset generation
1286 methodology incorporating human oversight, controlled data generation, and style transfer techniques.
1287 This appendix details our approach and its alignment with best practices in the literature.

1288 The `auto-STEER` methodology provides a systematic approach to generating datasets that mitigates
1289 the risk of data contamination, ensuring the integrity of benchmarks and the validity of results. Below,
1290 we outline the key aspects of `auto-STEER` that address this issue:

1291 **Challenging Models with Rephrasings:**

1292

1293 Rephrasings are known to cause significant variance in model performance, as demonstrated in the
1294 GSM-Symbolic dataset [Mirzadeh et al. \(2024\)](#) and other studies (e.g., [Zhu et al. \(2024\)](#), [Wang et al. \(2023\)](#)) highlighting how syntactic or stylistic changes can challenge generalization. In [Appendix G](#),
1295 we also show that much of the observed variance in LLM performance arises from these rephrasings,

1296 underscoring their role in robust evaluations. `auto-STEER` leverages this phenomenon to craft
1297 diverse rephrased questions that test beyond rote learning.

1298 **Systematic Question Generation:**

1299 `auto-STEER` generates new questions through a structured process that balances diversity and
1300 consistency. Questions are systematically rephrased or style-transferred to ensure they are different
1301 enough from the original templates to prevent memorization while retaining the same core meaning.
1302 This approach reduces the risk of overlap with pre-trained data while preserving the focus of the
1303 assessment.

1304
1305 The rapid advancement of large language models necessitates benchmarks that can evolve just
1306 as quickly. To address this, `auto-STEER` incorporates a user interface (UI) that allows users to
1307 regenerate entire datasets with minimal effort. By modifying domains, seeds, or even resampling
1308 numerical values, users can quickly produce a new dataset tailored to the latest needs. This adaptability
1309 ensures that benchmarks remain fresh and resistant to contamination as models advance.

1310 Through these features, `auto-STEER` provides a robust mechanism for creating datasets that challenge
1311 models in meaningful ways while maintaining a high degree of control over data integrity. The
1312 systematic generation of diverse, rephrased questions and the ability to regenerate datasets on demand
1313 make it a powerful tool for addressing data contamination in an era of rapidly evolving AI capabilities.

1314

1315 C TECHNICAL DESCRIPTIONS OF ADAPTATIONS

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

C.1 RPM (CONDITIONING):

Given the LLM’s output distribution over all possible tokens, filter to include only those that correspond to valid options. For example, if a question has four options then get the probabilities corresponding to ‘A’, ‘B’, ‘C’, and ‘D’. Then, compute softmax over the valid options to normalize the filtered probabilities into a distribution.

C.2 RPM (MIXING):

Alternatively, we restrict the output distribution to only valid option tokens O as follows: $\alpha \cdot p(o) + (1 - \alpha)^{1/|O|}$, where $o \in O$, $p(o)$ is the probability the LLM assigns to each token it outputs, and $\alpha = \sum_{o \in O} p(o)$. We then compute the softmax to normalize the resulting probabilities into a distribution.

In the mixing approach, if an LLM is confident in a valid option token the resulting distribution will place high probability on that token, but if an LLM places negligible probability on the valid option tokens then the resulting distribution will more closely resemble a uniform distribution.

C.3 PROMPT FOR GETTING ANSWERS FROM FREE-TEXT QA

User Message:

Report the answer that the following reasoning ended up with. Do not solve the question just look at the explanation text and report the answer if it exists or None if not. Just report the number or None.

Q: {question_text}

Explanation: {model_response}

D TECHNICAL DESCRIPTIONS OF METRICS

D.1 ACCURACY.

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

1350 D.1.1 EXACT-MATCH ACCURACY

1351
1352 This is the fraction of questions answered correctly.

1353 1354 D.1.2 NORMALIZED ACCURACY

1355 Elements can differ in their number of multiple choice options, leading to differences in the exact-
1356 match accuracy of random guessing. We can compensate for this by reporting the gap between the
1357 LLM’s exact-match accuracy and random guessing (Budesu & Bar-Hillel, 1993). We compute
1358 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
1359 the LLM outputs for question i , a_i is the indicator describing whether the top token is correct or not,
1360 N the number of questions in the element, and $|O_i|$ the number of options in the question. In other
1361 words, normalized accuracy rewards an LLM with 1 point for every correct answer and penalizes an
1362 LLM by 1 over the number of options minus 1 for each incorrect answer.

1363 1364 D.2 CALIBRATION

1365 It can also be useful to understand how confident an LLM is in its responses and the extent to which
1366 these confidence levels align with accuracy.

1367 1368 D.2.1 EXPECTED CALIBRATION ERROR

1369 We follow Liang et al. (2022) and Raman et al. (2024) in measuring the confidence of an LLM’s
1370 response and computing the expected calibration error (ECE; Naeini et al. 2015; Guo et al. 2017).
1371 ECE measures how closely the probability an LLM assigns to its top answer matches the actual
1372 probability of the correct answer, which in our case is 1. ECE first splits the data into M equally spaced
1373 bins, where each bin contains the probabilities the model assigned to their top token in that range:
1374 e.g., let p^{\max} be the set of most probable tokens for each question then if $M = 2$, then the first bin
1375 $B_1 = \{p \mid p \in p^{\max} \text{ and } p \in [0, 0.5]\}$. It is then defined as $\sum_{i \in [M]} |B_i|/N \cdot |\text{acc}(B_i) - \text{conf}(B_i)|$,
1376 where $\text{conf}(B_i)$ is the average probability the LLM assigned to its top token in bin B_i , and $\text{acc}(B_i)$
1377 denotes the exact-match accuracy in bin B_i . We allow users to choose the number of bins, however,
1378 we set $M = 10$ uniformly spaced over the interval $[0, 1]$ as is standard.

1379 1380 D.2.2 BRIER SCORE

1381 The Brier Score of an element is defined as

$$1382 \sum_{i=1}^N \frac{1}{|O_i|} \sum_{o \in O_i} (p_i(o) - a_i(o))^2,$$

1383 where $p_i(o)$ is the probability the LLM assigns to option o in question i . Thus, if an LLM is overly
1384 confident in an incorrect answer (e.g., assigns a probability of 0.9 to a wrong option), the Brier Score
1385 will penalize it more heavily.

1386 1387 D.2.3 EXPECTED PROBABILITY ASSIGNMENT

1388 EPA measures how much probability mass an LLM assigns to the correct answer option out of all
1389 possible options. It is defined as: $1/N \sum_{i \in [N]} p_i^*$, where p_i^* is the predicted probability that the LLM
1390 assigns to the correct option for question i .

1391 1392 D.3 ROBUSTNESS

1393 Elements are not the lowest level of granularity in our benchmark and aggregating on an element-by-
1394 element basis may hide where LLM performance varies. We introduce three robustness metrics that
1395 aggregate on the component fields of an element: domains, types, and perspectives. We compute the
1396 domain (type, perspective) robustness on each element by taking the minimum normalized accuracy
1397 over all domains (types, perspectives).

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

E EXAMPLES

E.1 GPT-4o 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-4o 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:

$$\text{Total Urban Demand} = 820 \times 5.416$$

$$\text{Total Urban Demand} = 4441.12$$

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:

$$\text{Total Suburban Demand} = 384 \times 2.528$$

$$\text{Total Suburban Demand} = 970.88$$

Step 3: Calculate the total demand

Now, sum the total demand from both markets:

$$\text{Total Demand} = \text{Total Urban Demand} + \text{Total Suburban Demand}$$

$$\text{Total Demand} = 4441.12 + 970.88$$

$$\text{Total Demand} = 5412$$

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

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

E.2 GPT-4o 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-4o 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.

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:

GPT-4o Message:

B

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

E.3 GPT-4O 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-4o 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.

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

$$\text{Expected Price} = 0.648 + 0.8988 + 1.1968 = 2.7436$$

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 ΔK that John might invest. The new capital K_2 will be $K_1 + \Delta K$.

The revenue function for a given capital K is:

$$\text{Revenue} = \text{Expected Price} \times \text{Production Output}$$

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

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:

$$\text{Cost} = (\Delta K)^2$$

Step 5: Calculate the Discounted Future Revenue

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

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

Step 6: Calculate the Net Profit

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

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

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

GPT-4o Message (Continued):

Step 7: Evaluate the Net Profit for Different Values of ΔK

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

Option A: $\Delta K = 0.44$

$$K_2 = 3.3 + 0.44 = 3.74$$

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

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

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

$$\text{Cost} = (0.44)^2 = 0.1936$$

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

Option B: $\Delta K = 0.14$

$$K_2 = 3.3 + 0.14 = 3.44$$

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

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

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

$$\text{Cost} = (0.14)^2 = 0.0196$$

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

Option C: $\Delta K = 1.1$

$$K_2 = 3.3 + 1.1 = 4.4$$

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

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

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

$$\text{Cost} = (1.1)^2 = 1.21$$

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

Option D: $\Delta K = 0.09$

$$K_2 = 3.3 + 0.09 = 3.39$$

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

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

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

$$\text{Cost} = (0.09)^2 = 0.0081$$

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

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

User Message:

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

A:

GPT-4o Message:

A

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

F FIGURES

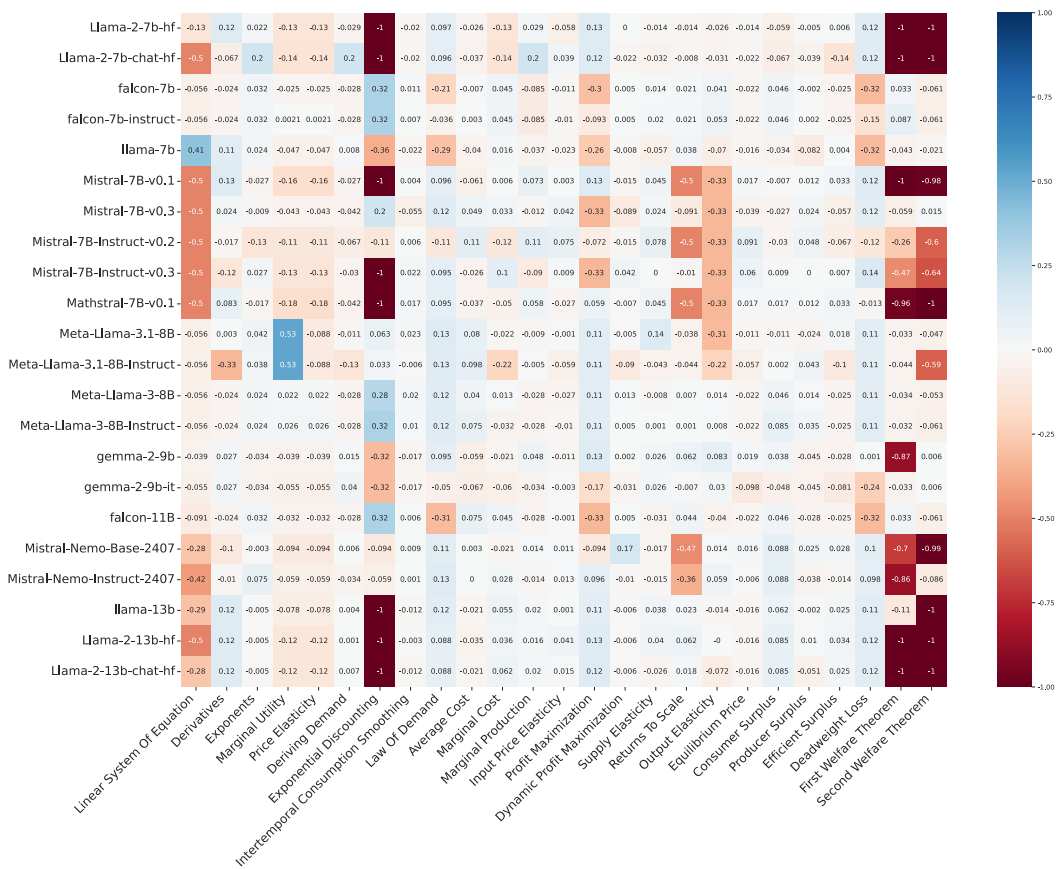


Figure 3: Heatmap of normalized accuracy of open-source models.

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

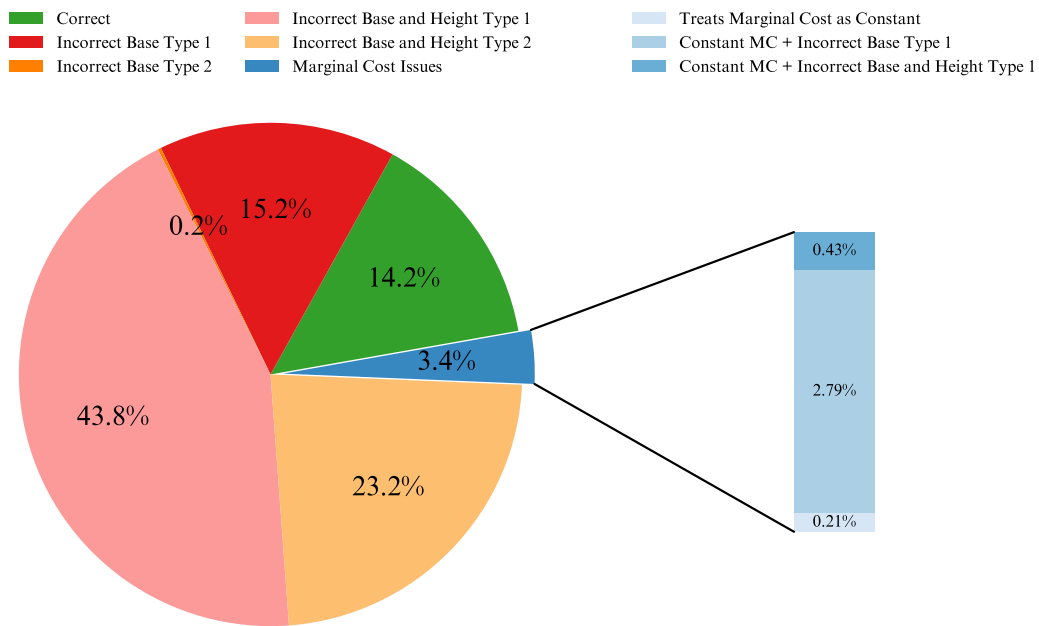


Figure 4: Error analysis of claude-3-5-sonnet on the Deadweight Loss of a Monopoly element. We further breakdown the errors by incorrectly interpreting the marginal cost.

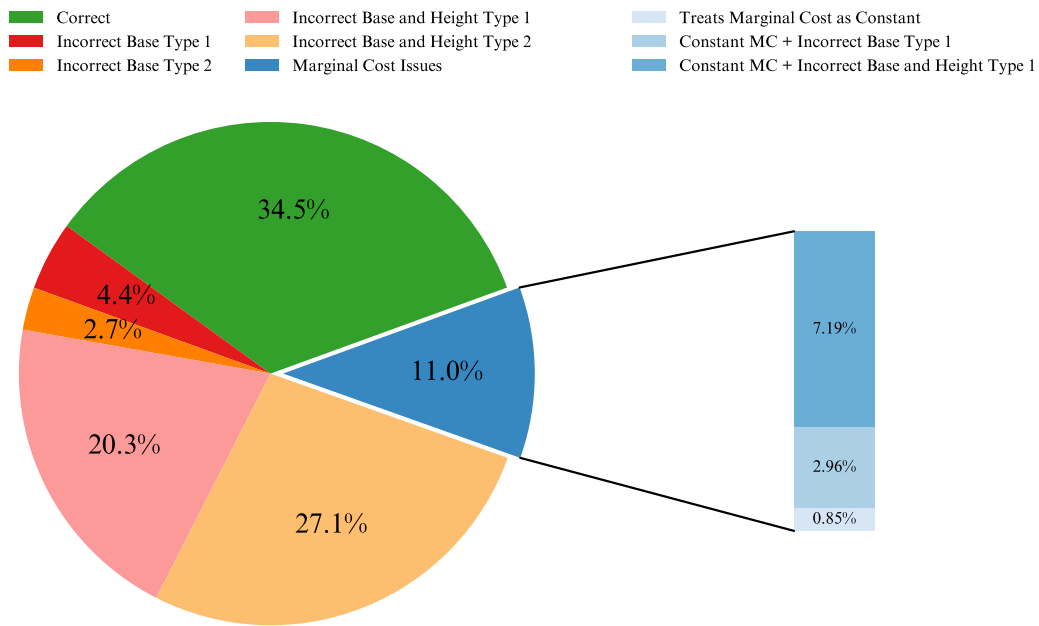


Figure 5: Error analysis of gpt-4o on the Deadweight Loss of a Monopoly element. We further breakdown the errors by incorrectly interpreting the marginal cost.

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

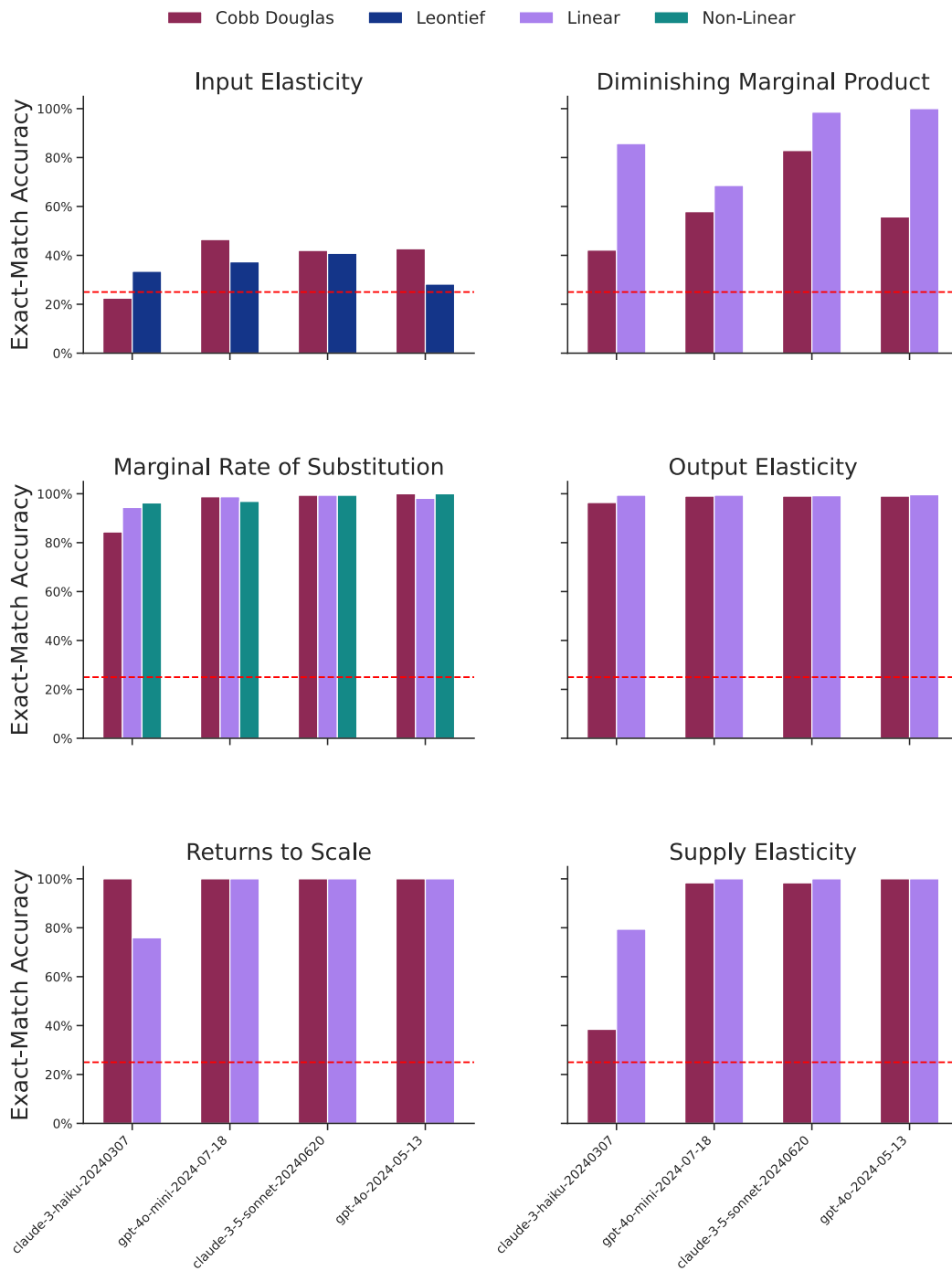


Figure 6: Exact-match accuracy of all closed-source models across six elements and four functional families (Cobb-Douglas, Leontief, Linear, and Non-Linear). The results demonstrate varying type robustness, with Cobb-Douglas being a generally more challenging functional family but not consistently harder 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.

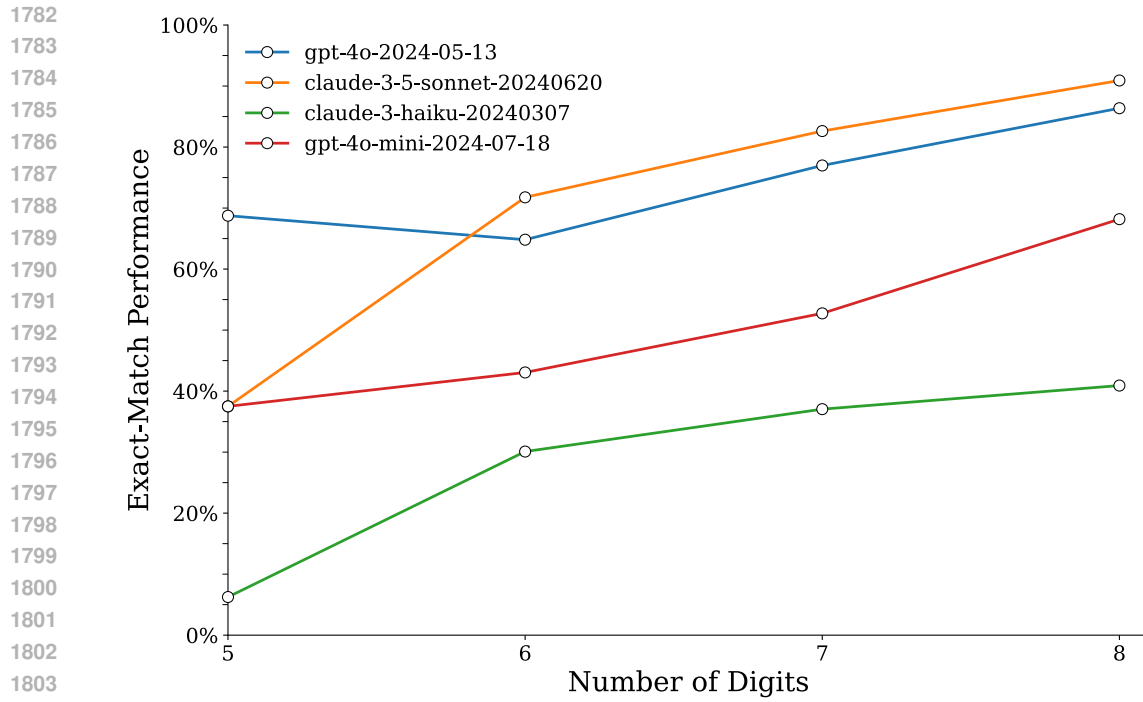


Figure 7: This figure depicts exact-match performance on the **Aggregation of Consumer Demand** element on the shown implementation of 0-CoT for the closed-source models against the number of digits of the correct answer.

G ANALYSIS OF REPHRASING VARIANCE

To understand the role of question rephrasings in our dataset, we conducted an analysis of variance (ANOVA) on all other controllable features. These features include type, domain, and perspective. The goal of this analysis was to quantify the variance in LLM performance attributable to these features and, by exclusion, infer the contribution of rephrasings to the remaining unexplained variance.

The results for the top-performing models, summarized in Table 2 through Table 6, indicate that the explained variance attributable to the controlled features is consistently low across all evaluated models. This leaves approximately 56% (for claude-3-5-sonnet) and up to 91% (for o1-preview) of the variance unexplained by the features included in the analysis. Given that question rephrasings are a systematic element of our dataset design and were not included as a feature in this analysis, we infer that the majority of this residual variance is due to differences in how models respond to semantically equivalent but syntactically varied prompts.

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		
		<i>R-squared</i>		0.0941
		<i>Adjusted R-squared</i>		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		
		<i>R-squared</i>		0.3368
		<i>Adjusted R-squared</i>		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		
		<i>R-squared</i>		0.2964
		<i>Adjusted R-squared</i>		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		
			<i>R-squared</i>	0.4436
			<i>Adjusted R-squared</i>	0.4430

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		
			<i>R-squared</i>	0.2336
			<i>Adjusted R-squared</i>	0.2330

Table 6: ANOVA Results for claude-3-haiku-20240307

H MODELS

Model Name	Model Card	Chat/ Instruction Tuned
Closed-Source		
<i>OpenAI</i>		
gpt-4o		✓
gpt-4o mini		✓
<i>Anthropic</i>		
claude-3-5-sonnet		✓
claude-3-haiku		✓
<i>Huggy Llama</i>		
llama-7b	huggyllama/llama-7b	×
llama-13b	huggyllama/llama-13b	×
llama-30b	huggyllama/llama-30b	×
llama-65b	huggyllama/llama-65b	×
<i>Meta Llama</i>		
Llama-2-7b-hf	meta-llama/Llama-2-7b-hf	×
Llama-2-13b-hf	meta-llama/Llama-2-13b-hf	×
Llama-2-7b-chat-hf	meta-llama/Llama-2-7b-chat-hf	✓
Llama-2-13b-chat-hf	meta-llama/Llama-2-13b-chat-hf	✓
Llama-3-8B	meta-llama/Meta-Llama-3-8B	×
Llama-3-8B-Instruct	meta-llama/Meta-Llama-3-8B-Instruct	✓
Llama-3.1-8B	meta-llama/Meta-Llama-3.1-8B	×
Llama-3.1-70B	meta-llama/Meta-Llama-3.1-70B	×
Llama-3.1-70B-Instruct	meta-llama/Meta-Llama-3.1-70B-Instruct	✓

Continued on next page

Model Name	Model Card	Chat/Instruction Tuned
<i>Mistral</i>		
Mistral-7B-v0.1	mistralai/Mistral-7B-v0.1	×
Mathstral-7B-v0.1	mistralai/Mathstral-7B-v0.1	×
Mistral-7B-v0.3	mistralai/Mistral-7B-v0.3	×
Mistral-7B-Instruct-v0.3	mistralai/Mistral-7B-Instruct-v0.3	✓
Mistral-Nemo-Base-2407 (12.2B)	mistralai/Mistral-Nemo-Base-2407	×
Mistral-Nemo-Instruct-2407 (12.2B)	mistralai/Mistral-Nemo-Instruct-2407	✓
<i>TIUAE</i>		
falcon-7B	tiiuae/falcon-7b	×
falcon-11B	tiiuae/falcon-11B	×
<i>AI21</i>		
Jamba-v0.1	ai21labs/Jamba-v0.1	×
AI21-Jamba-1.5-Mini	ai21labs/AI21-Jamba-1.5-Mini	×

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.

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

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-4o 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 8 shows the exact-match performance of both models on these three elements. We found no significant differences in performance between any of the models.

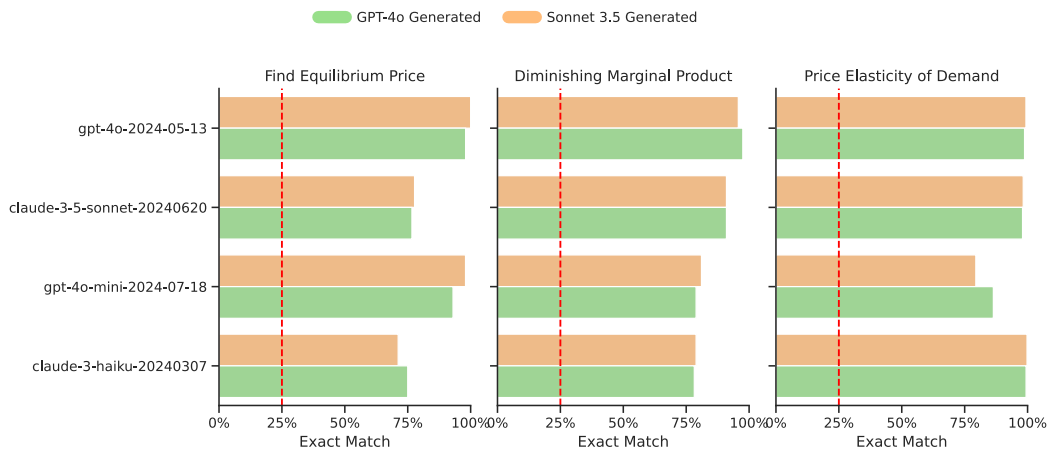


Figure 8: Exact-match performance comparison between closed-source models on three that were generated by gpt-4o and claude-3-5-sonnet. The elements were selected based on the observed performance differences across models, with **Price Elasticity of Demand** serving as a control. The figure shows no significant differences in performance between the models on these elements. Note that the red dotted line signifies random guessing performance.

I.2 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.

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

- **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.

I.3 PERFORMANCE GAPS BETWEEN FREE TEXT QA AND HIDDEN MCQA

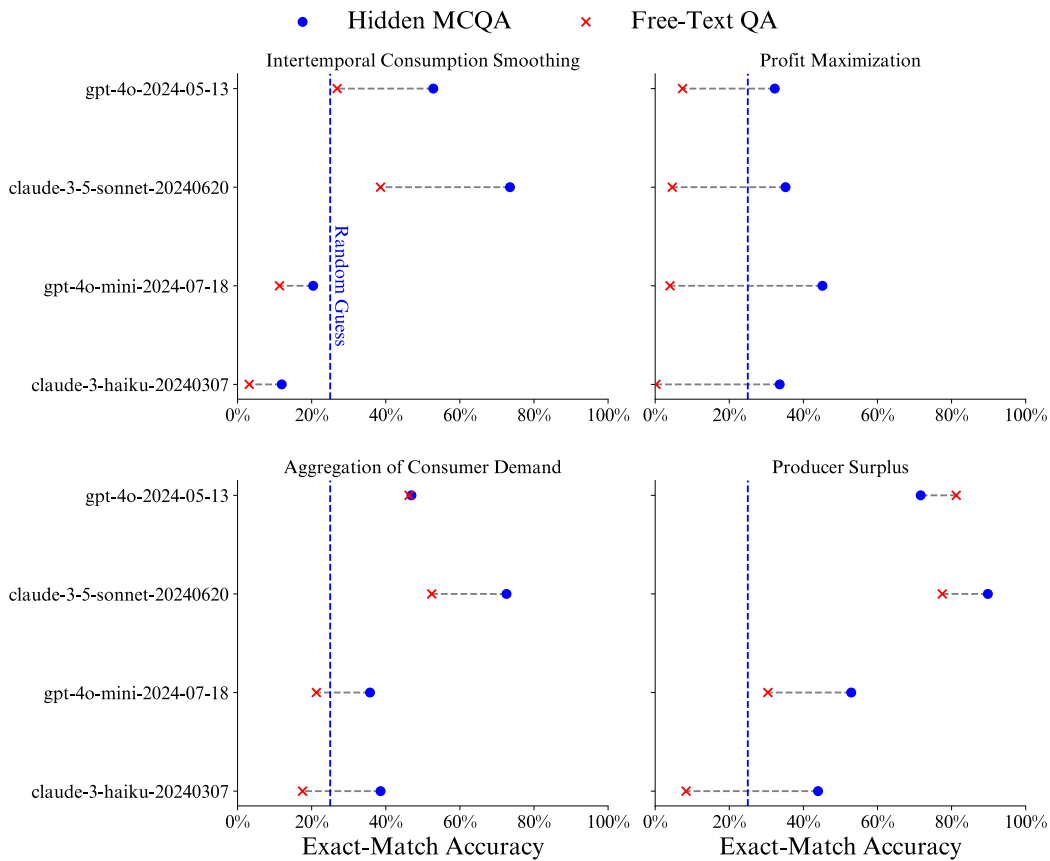


Figure 9: Comparison of exact-match accuracy for four elements (Intertemporal Consumption Smoothing, Profit Maximization, Aggregation of Consumer Demand, and Producer Surplus) across all closed-source models. The plot illustrates the difference in performance under the hidden and free-text adaptations, highlighting the impact of multiple-choice options on reasoning accuracy. Plotted in blue is the accuracy random guessing would have in the hidden adaptation, note that this line is not relevant for the free-text QA adaptation.

While most cases show a negative gap between hidden and free-text QA performance, there are notable exceptions. Figure 9 shows that in the Producer Surplus element, gpt-4o performed better in the free-text QA adaptation compared to the hidden adaptation. One might expect that the explanation for this positive gap is that the multiple-choice answers in the hidden adaptation were very similar, which may have caused confusion for the LLM. However, our analysis shows that when the free-text QA was scored correctly, gpt-4o selected an incorrect answer in the hidden adaptation with

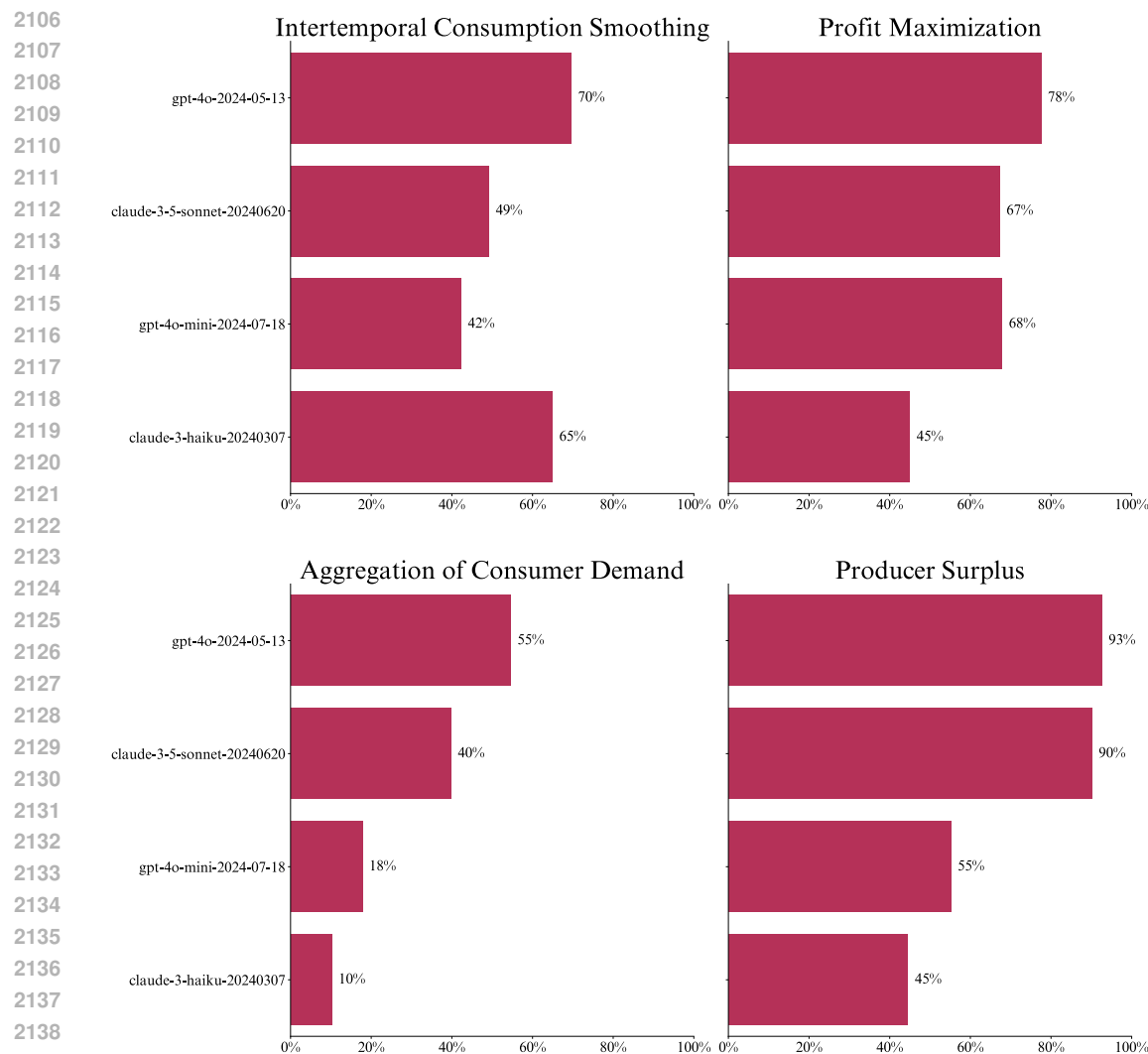


Figure 10: This figure depicts the percentage of time models were incorrect in the free-text adaptation but correct on the hidden adaptation due to choosing the closest answer. The plot compares the performance of four elements—Intertemporal Consumption Smoothing, Profit Maximization, Aggregation of Consumer Demand, and Producer Surplus—across all closed-source models.

at least a 10% difference from the correct answer 78% of the time. This suggests that while the LLM could reason effectively about the problem, it struggled to correctly match its reasoning to the multiple-choice options provided.

The plot also suggests that the performance gap in the Profit Maximization element is primarily due to the benefit random guessing has on accuracy in the hidden adaptation compared to free-text QA. Furthermore, in the Aggregation of Consumer Demand element, the inclusion of options after the reasoning step offered limited benefit, highlighting that the true advantage lies in including these options during the reasoning process.

These observations highlight an important nuance: although multiple-choice formats generally offer helpful structure for models, they may also hinder performance in certain scenarios.

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

J SECTION 3 IMAGES

User Prompt
Q: A baseball team is buying new equipment and needs baseballs, their demand for baseballs at any given price is expressed by the following demand function $-1.89Q + 2.6$. What is the team's consumer surplus if the price of baseballs is 1.24? Let's think step by step. Explain your reasoning.

Output
Model used: gpt-4
[Detailed reasoning steps]

User Prompt
A. 0.49
B. 0.33
C. 0.21
D. 1.34
Answer by writing the option letter corresponding to the correct option. WRITE ONLY A SINGLE LETTER.
A:

Output
Model used: gpt-4
D

User Prompt
Q: A baseball team is buying new equipment and needs baseballs, their demand for baseballs at any given price is expressed by the following demand function $-1.89Q + 2.6$. What is the team's consumer surplus if the price of baseballs is 1.24?
A. 0.49
B. 0.33
C. 0.21
D. 1.34
Let's think step by step. Explain your reasoning.

Output
Model used: gpt-4
[Detailed reasoning steps]

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

Output
Model used: gpt-4
B

Figure 11: (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.

2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267

K WEB APPLICATION

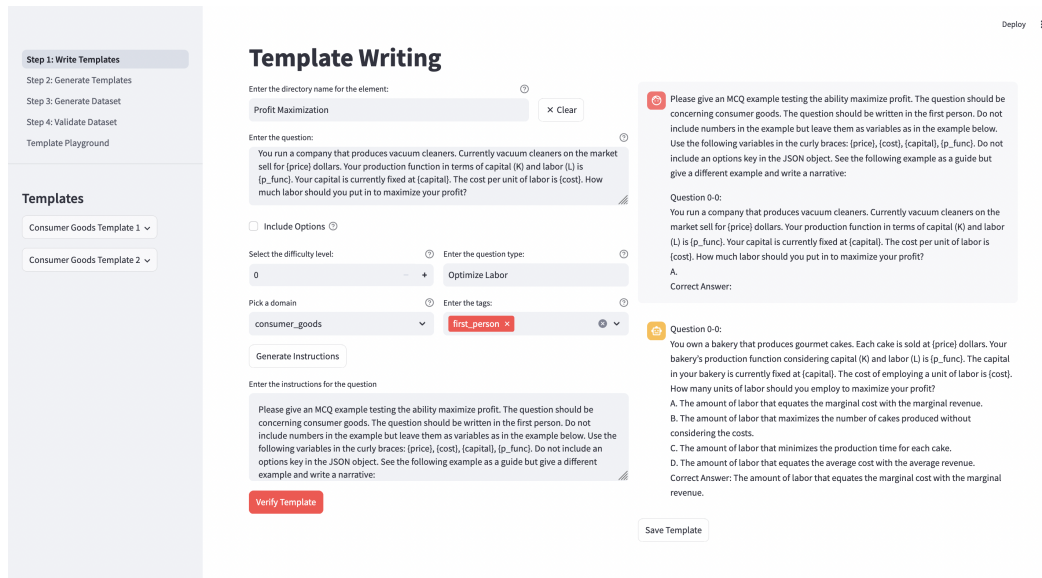


Figure 12: The web app user interface for template writing. This page includes fields for type, domain, grade level and tags (including perspectives). The right shows an example of template verification which uses a LLM to generate another template using the example seed.

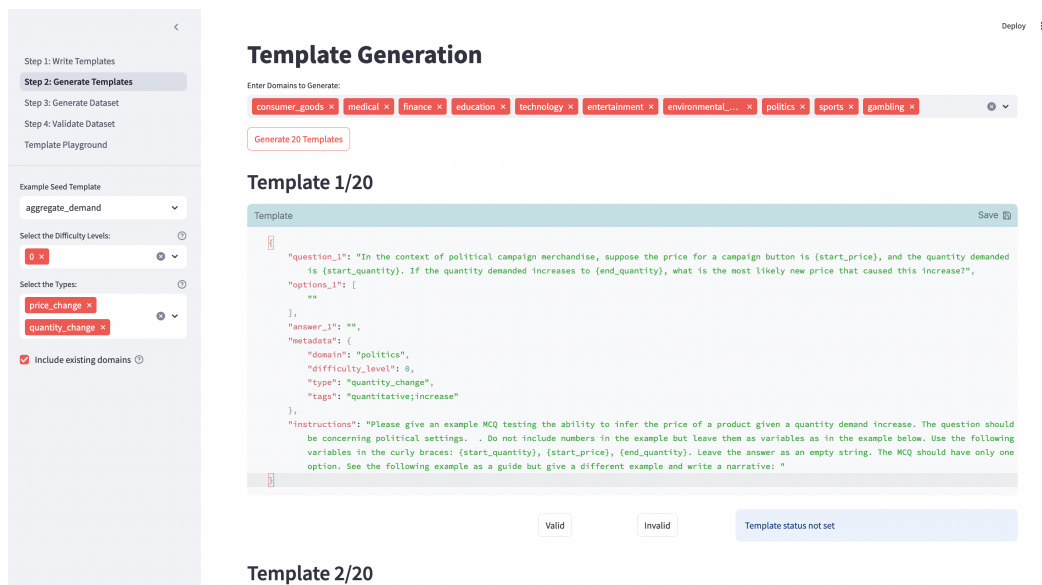


Figure 13: The web app user interface for template generation. This page allows for a selection of domains, and types for which templates will be generated using the available example seeds. Templates can then be verified and saved by the user.