

交渉 LLM における ZOPA 中間点での合意形成のための プロンプト・エンジニアリング

園田哲也¹ 宇津呂武仁² 鈴木良弥¹

¹山梨大学 工学部 メカトロニクス工学科 ²筑波大学 システム情報系 知能機能工学科

概要

大規模言語モデル (LLM) による交渉では、交渉可能領域 (ZOPA) を十分に考慮できず、自身の限界価格に極端に固着してしまう傾向が指摘されている。本論文では、例示により双方の譲歩によって ZOPA 中間点での合意形成を促すプロンプト・エンジニアリング手法の有効性を検証した。実際の人間の交渉を例示する few-shot 手法では、同一価格帯では有効だが、異なる価格帯では効果が限定的であった。例示から傾向を抽出し記述させることで、状況に合わせて目標合意価格を導出させる Chain-of-Thought (CoT) 手法を導入した結果、価格帯や ZOPA が異なる設定においても、ZOPA 中間点での合意形成能力が向上することを確認した。

1 はじめに

大規模言語モデル (LLM) は、人間の交渉トレーニングや交渉の支援エージェントとしての活用が期待されている。こうした人間との協調を前提とした用途においては、人間的な傾向を掴んだ交渉を行う能力が不可欠である。

しかし Shah ら [10] は、住宅売買シナリオにおいて人間は交渉可能領域 (ZOPA: Zone of Possible Agreement) 中間点で公平に合意する一方、LLM は自身の限界価格に極端に固着し、柔軟な譲歩が行えない傾向があることを示した。本論文では、プロンプト・エンジニアリングによりこの問題を解決し、LLM エージェントによる ZOPA 中間点での合意形成を実現することを目的とする。

2 関連研究

LLM を用いた交渉エージェントに関しては、その行動特性や課題に関する分析が多数報告されている [9, 1, 7, 4]。代表的な例として、Ross ら [9] は効用理論に基づき LLM のバイアスを調査し、不公平・

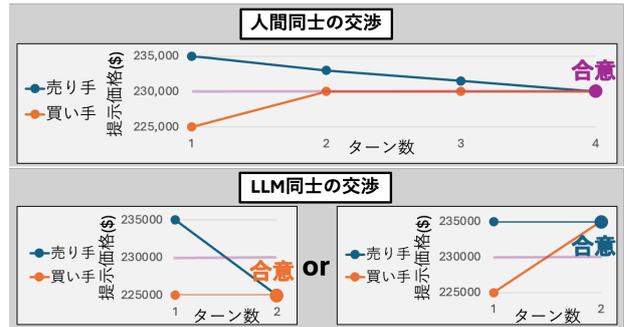


図 1: 交渉タスクにおける人間と LLM の挙動の違い

リスク・損失回避等の観点で、その行動は経済合理的とも人間的とも言えないと報告している。

こうした課題に対し、戦略的な学習や推論により交渉能力を向上させる試みも活発に行われている [3, 12, 6, 8, 2, 11]。例えば、Fu ら [3] は、自己対戦と AI のフィードバックを用いた文脈内学習により、対話ごとに戦略を改善し、より有利な価格での合意を目指す手法を提案しており、武並ら [11] は、目標額よりも意図的に高い初期提示額を設定して相手に提示することで、交渉結果を自身に有利に導くことを検証した。

これらは、LLM が交渉において可能な限り自身の利益を最大化することを目指すアプローチである。

一方、Shah ら [10] はこれらとは対照的に、譲歩ダイナミクスをモデル化する数理的枠組みと定量的指標を導入した人間との比較実験を通じて、LLM が初期条件に縛られ、人間のような柔軟な譲歩ができない問題 (アンカリング) を指摘している。本論文では、自身の利益を最大化するのではなく、このアンカリング問題を解消し、ZOPA 中間点での合意を実現するためのプロンプト手法を提案する。

3 価格交渉タスク

本論文では、Shah ら [10] が LLM のアンカリング問題を実証した、Heddaya ら [5] の住宅価格交渉タスクを採用する。本タスクは同条件での人間の

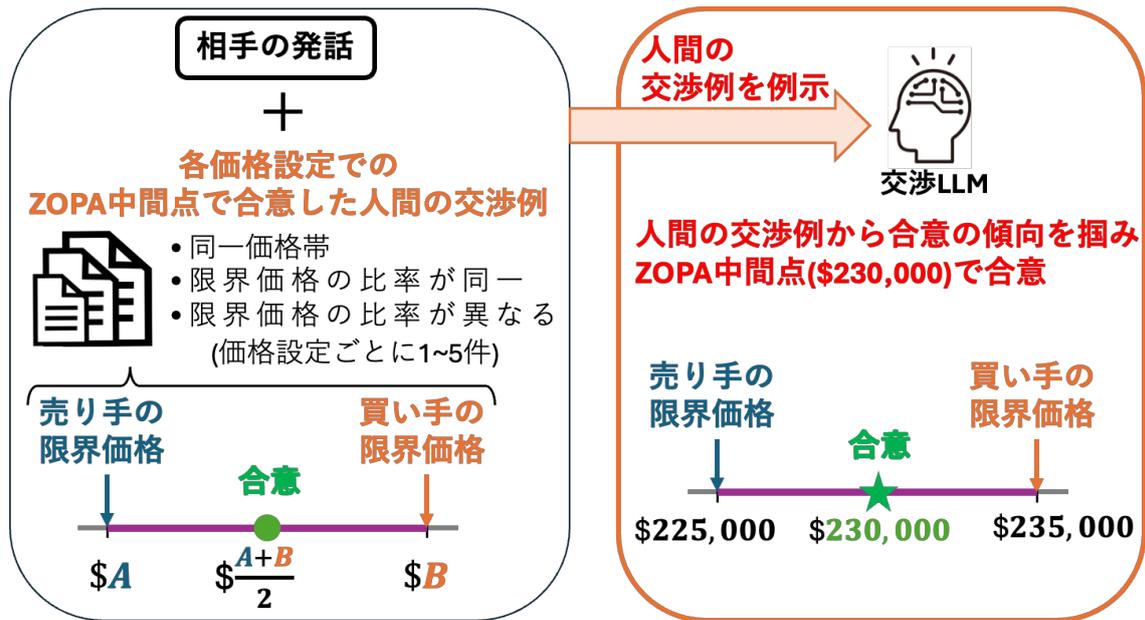


図 2: 本論文の目的: 異なる価格設定の交渉例の例示により ZOPA 中間点での合意を目指す

表 1: few-shot プロンプトにおける 5 つの価格設定
(L_S : 売り手限界価格, L_B : 買い手限界価格, ZOPA: [L_S, L_B])

価格設定	概要	ZOPA
同一価格帯	全ての例で ZOPA を変更しない.	[\$225,000, \$235,000]
同一比率	L_S, L_B を同じ倍率で変更.	
同一 ZOPA	全ての例で共通の倍率を使用.	[\$900,000, \$940,000]
異なる ZOPA	例ごとに, 異なる倍率を使用.	[\$450,000, \$470,000] や [\$675,000, \$705,000] など
異なる比率	L_S, L_B それぞれ異なる倍率で変更.	
同一 ZOPA	全ての例で共通の変更を使用.	[\$180,000, \$1,245,000]
異なる ZOPA	例ごとに, 異なる倍率の組み合わせを使用.	[\$98, \$105] や [\$76,900, \$86,000] など

交渉データが利用可能であり, 人間と LLM の挙動比較に適している. Heddaya らの設定では, 売り手の許容範囲が \$225,000 から \$240,000, 買い手の許容範囲が \$225,000 から \$235,000 と条件が非対称だが, 実質的な合意可能領域 (ZOPA) は \$225,000 から \$235,000 の間であった. 本論文では, この ZOPA 範囲を維持しつつ, 売り手と買い手の価格条件を共に \$225,000 から \$235,000 とすることで対称化を行った. これは, 価格帯の非対称性による攪乱要因を排除し, 純粋なアンカリング問題の影響を検証するためである.

図 1 に本タスクにおける典型的な挙動を示す. 人間は互いに歩み寄り ZOPA 中間点付近で合意する傾向があるのに対し, LLM(zero-shot) はこの対称な状況下においても自身の限界価格に固着し, 極端な価格で合意してしまうことが確認された. 本論文の目的は, 図 2 に示すように, 例示を通じてこの固着

を解消し, 異なる価格設定においても人間と同様の ZOPA 中間点での合意を実現することである.

4 ZOPA 中間点での合意形成のためのプロンプト・エンジニアリング

本章では, 前章で述べた LLM のアンカリング問題を解決し, ZOPA 中間点での合意形成を実現するためのプロンプト設計について述べる.

4.1 few-shot プロンプティング

"Language of Bargaining" データセット [5]¹⁾ から, ZOPA 中間点である \$230,000 で合意に至った対話データ 5 件を選定し, LLM にプロンプトとして提示する. 本論文では, 選定した対話データの文脈を維持しつつ, 価格設定のみを人手により変更した 5 パターンのデータセットを作成した (表 1). (詳細は付

1) <https://huggingface.co/datasets/ChicagoHAI/language-of-bargaining>

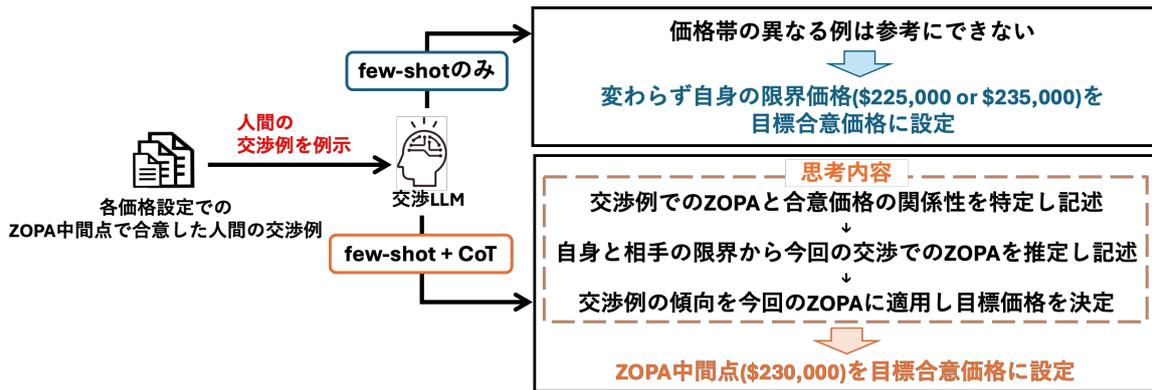


図 3: 異なる価格設定の交渉例の有効性: few-shot のみでは効果なし・few-shot+CoT では有効

録の表 2 参照).

4.2 思考連鎖

few-shot のみでは、同一価格帯の設定ではある程度改善するものの、後述する実験結果が示すように、価格設定が異なる例を現在の交渉の参考にすることができない。そこで、思考の過程を明示的に記述させることで、適切な目標合意価格に到達させるプロンプト(付録 A を参照)を構築した。

具体的には、交渉 LLM の発話前に図 3 に示す以下のプロセスを思考させた。

1. 交渉例の ZOPA と合意価格の関係を特定し記述
提示された few-shot 事例について、交渉が公平に行われたと仮定して相手の限界価格を推測し、自身の限界価格と合わせて、事例ごとの ZOPA(合意可能領域)と、ZOPA 幅に対する利益率を特定し記述させる。
2. 今回の交渉における ZOPA を推定し記述
現在の自身の限界価格と、対話履歴および売り出し価格から推定される相手の限界価格を明記させ、今回の交渉における ZOPA を定義させる。
3. 目標合意価格と次の発話を決定
先に特定した交渉例の利益率を現在の ZOPA に適用し、具体的な目標合意価格を算出させる。

5 評価

OpenAI 社の GPT-4.1²⁾を用いた自己対戦により、前章の 5 設定かつ shot 数 (1-5) の組み合わせで各 20 回試行した。temperature=1.0, 最大ターン数 10 に設定し、10 ターン以内に合意できなかった場合は

2) <https://platform.openai.com/docs/models/gpt-4.1>

交渉決裂とした。比較手法は zero-shot, few-shot のみ, few-shot+CoT とし、評価指標には ZOPA 中間点での合意率および二乗平均平方根誤差 (RMSE: Root Mean Squared Error) を採用した。ここで ZOPA 中間点とは、双方の限界価格 (\$225,000 および \$235,000) の中央値である \$230,000 を指す。

目標合意価格 $y_{target} = 230,000$, 合意価格の平均値 \bar{y} , 合意価格の標準偏差 σ とした際、RMSE は以下の式で算出される。

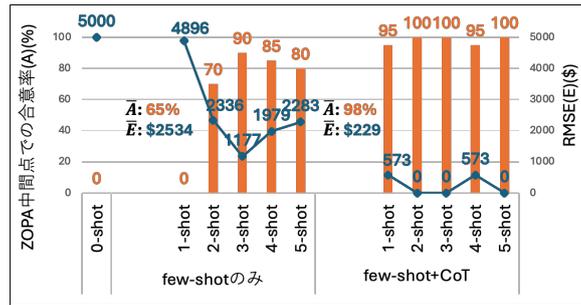
$$RMSE = \sqrt{(\bar{y} - y_{target})^2 + \sigma^2} \quad (1)$$

この値が小さいほど、アンカリング問題による偏りやばらつきが解消されたことを示す。

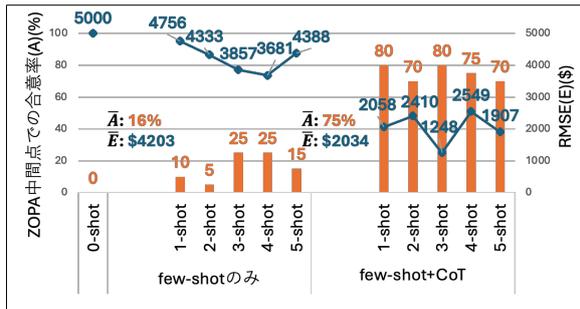
実験結果を図 4 に示す。(詳細は付録の表 3 参照)。Shah ら [10] の報告と同様に、zero-shot では ZOPA 中間点での合意はなく、売り手の限界価格での合意のみだった。また、few-shot のみでは、同一価格帯の設定において合意率 65%, RMSE\$2,534 と一定の改善が見られたが、価格帯や比率を変更したその他の設定では、合意率は 2% - 16% となり、RMSE も \$4,000 以上と大きな改善は見られなかった。この結果は、few-shot のみでは、例での比率や ZOPA の構造を、設定の異なる自身の状況へ転用できていないことを示している。

これに対し、few-shot+CoT を適用した結果、全ての設定において顕著な性能向上が確認された。特に、「異なる比率・異なる ZOPA」の設定においては、平均合意率が few-shot のみの 2% から 83% へと大幅に向上し、5-shot の場合には合意率 100%, RMSE\$0 を達成するなど、十分な例示があれば異なる価格設定であっても構造を把握し、最適な合意形成が可能であることが実証された。

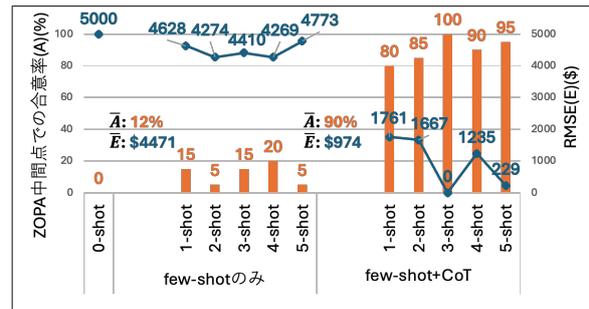
また、例示データの構成に関しても、興味深い傾



(a) 同一価格帯

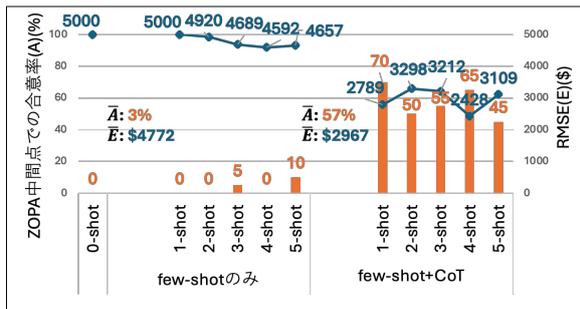


(b-1) few-shot間でZOPAを統一

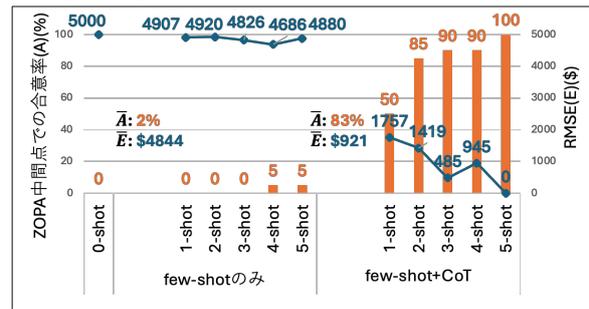


(b-2) few-shot間でZOPAを統一しない

(b) 売り手・買い手の限界価格の比率が同一



(c-1) few-shot間でZOPAを統一



(c-2) few-shot間でZOPAを統一しない

(c) 売り手・買い手の限界価格の比率が異なる

図4: 交渉例の5種価格設定における zero-shot, few-shot のみ, few-shot+CoT の比較 (A: ZOPA 中間点での合意率 (%), E: RMSE(\$))

向が確認された。全ての例で同一の ZOPA を用いた場合と比較して、例ごとに異なる ZOPA を用いた場合の平均合意率は「同一比率」の設定で 75% から 90% へ、「異なる比率」の設定でも 57% から 83% へと大幅な改善が見られた。通常、学習データはテスト条件に近い方が有利とされるが、本タスクにおいては、提示される数値にバラつきがあることで、モデルの特定の数値への表面的な模倣を防ぎ、より抽象度の高い構造の解析が促進されたものと推察される。

6 おわりに

本論文では、交渉タスクを行う LLM が自身の限界価格に固着するアンカリング問題を解決するため、ZOPA 中間点での合意を目指すプロンプト手法

を検証した。

実験の結果、few-shot のみでは、例示と異なる価格設定への適応が困難であることが確認された。これに対し、few-shot+CoT では、交渉の文脈を構造的な情報に変換して推論させることで、価格帯や比率が異なる設定においても、ZOPA 中間点での合意が可能であることを実証した。また、例示データの構成に関する分析から、同一 ZOPA の例示を繰り返すよりも、あえて異なる ZOPA の例を提示した方が高い汎化性能を示すことが明らかとなった。

本論文では中間点での合意事例のみを用いたが、今後は異なる妥協点で合意した事例を用いた場合の影響や、双方の ZOPA が重ならない困難な条件下での挙動検証などを計画している。

参考文献

- [1] Federico Bianchi, Patrick John Chia, Mert Yuksekgonul, Jacopo Tagliabue, Dan Jurafsky, and James Zou. How well can LLMs negotiate? negotiation arena platform and analysis. In *Proc. 41st ICML*, pp. 3935–3951, 2024.
- [2] Jiangjie Chen, Siyu Liu, Xuhui Zhu, and Yunjia Yao. Put your money where your mouth is: Evaluating strategic planning and execution of LLM agents in an auction arena. In *Proc. ACL 2024 Findings*, pp. 4215–4233, 2024.
- [3] Yao Fu, Haoran Peng, Tushar Khot, and Mirella Lapata. Improving language model negotiation with self-play and in-context learning from AI feedback. *arXiv preprint arXiv:2305.10142*, pp. 1–13, 2023.
- [4] Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah D. Goodman. Understanding social reasoning in LLMs with the ultimatum game. In *Proc. 45th CogSci*, pp. 1–14, 2023.
- [5] Mourad Heddaya, Solomon Dworkin, Chenhao Tan, Rob Voigt, and Alexander Zentefis. Language of bargaining. In *Proc. 61st ACL*, pp. 13161–13185, 2023.
- [6] Yunbo Long, Liming Xu, Lukas Beckenbauer, Yuhan Liu, and Alexandra Brintrup. EvoEmo: Towards evolved emotional policies for LLM agents in multi-turn negotiation. *arXiv preprint arXiv:2509.04310*, pp. 1–19, 2025.
- [7] Nicolas Lorè and Babak Heydari. Strategic behavior of large language models: Game structure vs. contextual framing. *arXiv preprint arXiv:2309.05898*, pp. 1–16, 2023.
- [8] Edoardo De Luigi, Ildiko Pilan, Danilo Giampiccolo, and Marco Guida. NegotiationToM: A benchmark for theory of mind in negotiation with LLMs. In *Proc. LREC-COLING 2024*, pp. 11822–11832, 2024.
- [9] Jillian Ross, Yoon Kim, and Andrew W. Lo. LLM economicus? mapping the behavioral biases of LLMs via utility theory. In *Proc. 1st COLM*, pp. 1–22, 2024.
- [10] Cheril Shah, Akshit Agarwal, Kanak Garg, and Mourad Heddaya. LLM rationalis? measuring bargaining capabilities of AI negotiators. In *Proc. NeurIPS 2025 Workshop on Multi-Turn Interactions in LLMs*, pp. 1–10, 2025.
- [11] 武並佳輝, Yin Jou Huang, 村脇有吾, Chenhui Chu. LLMによる価格交渉シミュレーションにおけるアンカリング効果の検証. 言語処理学会 第 31 回年次大会論文集, pp. 921–925, 2025.
- [12] Tian Xia, Zhiwei He, Tong Ren, Yibo Miao, Zhuosheng Zhang, Yang Yang, and Rui Wang. Measuring bargaining abilities of LLMs: A benchmark and a buyer-enhancement method. In *Proc. ACL 2024 Findings*, pp. 3579–3602, 2024.

A 使用したプロンプト

売り手・買い手、周辺情報などのプロンプトは、Shah ら [10]³⁾ の記述を元に作成した。

以下に、本論文で導入した CoT プロンプトの詳細を示す。なお、紙面の都合上、記述は一部要約している。

交渉戦略を立てるために、以下の順序で思考を記述せよ。

【Step 0: 成功事例の構造解析】

[事例のデータ抽出と利益計算] 提供された全事例について、(A) 自分の限界と (B) 最終合意価格を抽出し、(E) 自分の確定利益額 (買い手: $A - B$, 売り手: $B - A$) を計算せよ。

[相手の限界価格の逆算] 「交渉は公平な妥協点で合意する」という前提から、相手にも自分と同じ程度の利益 (E) があると仮定し、(C) 相手の推定限界価格を逆算せよ

[成功事例での ZOPA] (A) 自分の限界、(C) 相手の推定限界から、(D) ZOPA (合意可能範囲) を定義せよ。

[自分の利益率の計算] ZOPA 全体の幅のうち「自分がどれだけの利益を確保できたか」を計算し、% で記述せよ (式: $\text{利益率}(\%) = \text{利益額}(E) \div \text{ZOPA 幅}(|A - C|)$)。

【Step 1: 今回の ZOPA の推定】

[相手限界の推定] Step 0 の構造を今回の設定に当てはめ、今回の相手の限界価格を記述せよ。相手の提示額ではなく「売り出し価格」等に連動する真の限界 (予算/底値) を見抜け。

[ZOPA の特定] 自分の限界価格と今回の相手の限界価格から、今回の ZOPA を定義せよ。

【Step 2: 目標合意価格 (Target) の設定】

[Target の算出] Step 0 の「自分の利益率 (%)」を今回の ZOPA に当てはめ、Target を算出せよ (例: 買い手なら 上限 - (ZOPA 幅 × 利益率))。

【Step 3: 次の行動の決定】

[ギャップ分析] 相手の最新提示額と Target の差を確認せよ。

[安易な妥協の回避] 相手の提示が「自分の限界」に近いなら拒否し、Target と一致するまで、論理的に価格修正を求めよ。

[次の提示価格] Target へ誘導するための価格を提示せよ。相手の提示額が Target より不利な場合のみ、Target 価格を提示して粘り強く交渉せよ。

B 価格設定・実験結果

すべての例の価格設定を表 2 に、実験結果を表 3 に示す。

表 2: few-shot 例の価格設定

設定 / 例	売手限界 (\$)	買手限界 (\$)	合意価格 (\$)
同一価格帯			
全例	225,000	235,000	230,000
同一比率・同一 ZOPA			
例 1-5	900,000	940,000	920,000
同一比率・異なる ZOPA			
例 1	450,000	470,000	460,000
例 2	675,000	705,000	690,000
例 3	900,000	940,000	920,000
例 4	1,125,000	1,175,000	1,150,000
例 5	1,350,000	1,410,000	1,380,000
異なる比率・同一 ZOPA			
例 1-5	180,000	1,245,500	712,750
異なる比率・異なる ZOPA			
例 1	98	105	101.5
例 2	3,210	4,000	3,605
例 3	76,900	86,000	81,450
例 4	553,000	581,000	567,000
例 5	6,380,000	8,600,000	7,490,000

表 3: 各設定・手法・例示数ごとの詳細結果

設定	手法	shot 数	合意率 (%)	RMSE (\$)
Baseline	zero-shot	0	0	5000.00
		1	0	4896.00
		2	70	2336.00
		3	90	1177.00
		4	85	1979.00
同一価格帯	few-shot のみ	5	80	2283.00
		1	95	573.00
		2	100	0.00
		3	100	0.00
		4	95	573.00
	few-shot+CoT	5	100	0.00
		1	10	4756.00
		2	5	4333.00
		3	25	3857.00
		4	25	3681.00
同一比率 同一 ZOPA	few-shot のみ	5	15	4388.00
		1	80	2058.00
		2	70	2410.00
		3	80	1248.00
		4	75	2549.00
	few-shot+CoT	5	70	1907.00
		1	15	4628.00
		2	5	4274.00
		3	15	4410.00
		4	20	4269.00
同一比率 異なる ZOPA	few-shot のみ	5	5	4773.00
		1	80	1761.00
		2	85	1667.00
		3	100	0.00
		4	90	1235.00
	few-shot+CoT	5	95	229.00
		1	0	5000.00
		2	0	4920.00
		3	5	4689.00
		4	0	4592.00
異なる比率 同一 ZOPA	few-shot のみ	5	10	4657.00
		1	70	2789.00
		2	50	3298.00
		3	55	3212.00
		4	65	2428.00
	few-shot+CoT	5	45	3109.00
		1	0	4907.00
		2	0	4920.00
		3	0	4826.00
		4	5	4686.00
異なる比率 異なる ZOPA	few-shot のみ	5	5	4880.00
		1	50	1757.00
		2	85	1419.00
		3	90	485.00
		4	90	945.00
	few-shot+CoT	5	100	0.00

3) <https://arxiv.org/abs/2512.13063>