

From Text to Trust: Empowering AI-assisted Decision Making with Adaptive LLM-powered Analysis (Supplementary Material)

Zhuoyan Li
li4178@purdue.edu
Purdue University
West Lafayette, Indiana, USA

Hangxiao Zhu
hangxiao@tamu.edu
Texas A&M University
College Station, Texas, USA

Zhuoran Lu
lu800@purdue.edu
Purdue University
West Lafayette, Indiana, USA

Ziang Xiao
ziang.xiao@jhu.edu
Johns Hopkins University
Baltimore, Maryland, USA

Ming Yin
mingyin@purdue.edu
Purdue University
West Lafayette, Indiana, USA

ACM Reference Format:

Zhuoyan Li, Hangxiao Zhu, Zhuoran Lu, Ziang Xiao, and Ming Yin. 2025. From Text to Trust: Empowering AI-assisted Decision Making with Adaptive LLM-powered Analysis (Supplementary Material). In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*, April 26–May 1, 2025, Yokohama, Japan. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3613904.3642625>

A EMPIRICAL EXAMINATIONS OF IMPACTS OF LLM-POWERED ANALYSIS IN AI-ASSISTED DECISION MAKING (ADDITIONAL DETAILS)

The full demographic information and statistics of participants in the study of this phase (where we have three treatments, i.e., CONTROL, SEQ, and ALL) for the income prediction and recidivism prediction tasks are shown in Table A.1 and Table A.2, respectively.

B COMBINING HUMAN DECISIONS AND AI PREDICTIONS

In the main paper, we aim to measure the reliability of AI model recommendations. Following previous research [2, 6], we combined independent human decisions with AI model predictions to determine the targeted decision for each task instance. We evaluated the human-AI combination method [2] and several truth inference methods used in crowdsourcing for truth discovery. We detailed the process of evaluation below.

Simulating Human Independent Decision. To understand how humans independently make decisions on the task instance, we first conducted a pilot study on Prolific to collect independent human decision behavior data on income prediction tasks and recidivism prediction tasks. We recruited 40 participants for each task. Each recruited participant needed to complete 15 tasks. With the collection of human behavior data, we then fitted two-layer neural networks to simulate human independent decision behavior. We optimized

these independent behavior models using Adam [3] with an initial learning rate of $1e-4$ and a batch size of each training iteration of 128. The number of training epochs is set as 10. The average accuracy of 5-fold validation for the fitted decision models is 0.81 for the income prediction task and 0.84 for the recidivism task, both of which we found to be satisfactory. We then utilized these fitted models to simulate independent human decisions $y_{\text{independent}}^h$ in the human-AI combination process to determine the potentially better decisions.

Comparing Combination Performance. We consider the human + AI combination method proposed in [2] and a few truth inference methods in crowdsourcing as baselines in the evaluation, including GLAD [7], CATD [4], LFC [5], EM [8], and MV [8]. These methods combine the human independent decisions $y_{\text{independent}}^h$ predicted by the fitted independent human behavior models and AI model recommendations y^m to produce combined decisions y_{combine} . The accuracy of y_{combine} on holdout task pools when using different methods to combine humans' (predicted) independent decisions and AI's decision recommendation is reported in Table B.1. In general, we found that the human + AI combination method proposed in [2] outperforms other baselines. By integrating human decisions with AI predictions, this method shows superior performance to either AI solo or human solo across the two types of decision making tasks. Consequently, we used the combined decisions y_{combine} from the human + AI combination method as the targeted decision in subsequent experiments to select the LLM rationale analysis.

C ALGORITHMIC SELECTION OF LLM-POWERED ANALYSIS IN AI-ASSISTED DECISION MAKING (ADDITIONAL DETAILS)

C.1 Algorithmic Framework Setting

The initial hidden state distribution $\mathcal{P}(\mathbf{h}_0|\mathbf{x}, y^m; \theta_{\text{init}})$ is modeled as a Gaussian distribution conditioned on the task \mathbf{x} and the AI model prediction y^m . Specifically, it is represented as $\mathcal{N}(\mathbf{h}^0; \mu_{\theta_{\text{init}}}(\mathbf{x}, y^m), \Sigma_{\theta_{\text{init}}}(\mathbf{x}, y^m))$, where $\mu_{\theta_{\text{init}}}$ and $\Sigma_{\theta_{\text{init}}}$ are parameterized by one-layer feedforward networks. For state updating, we also use a Gaussian distribution to characterize this process as $\mathcal{N}(\mathbf{h}^t; \mu_{\theta_{\text{update}}}(\mathbf{r}^t, \mathbf{a}^t, \mathbf{h}^{t-1}), \Sigma_{\theta_{\text{update}}}(\mathbf{r}^t, \mathbf{a}^t, \mathbf{h}^{t-1}))$. Here, \mathbf{r}^t is encoded as three parts: (1) the text embedding of the LLM-powered analysis encoded with a BERT model [1]; (2) the task

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
CHI '25, April 26–May 1, 2025, Yokohama, Japan
© 2025 Copyright held by the owner/author(s).
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM.
<https://doi.org/10.1145/3613904.3642625>

Demographics		CONTROL (<i>N</i> = 41)	SEQ (<i>N</i> = 45)	ALL (<i>N</i> = 48)
Gender	Male	26.9%	53.3%	37.5%
	Female	73.1%	44.4%	60.4%
	Others	0	2.3%	2.1%
Age	Below 35	60.9%	46.7%	54.1%
	35 -44	29.2%	33.3%	16.7%
	45 or above	9.9%	20%	29.2%
Race	White	51.2%	68.8%	58.4%
	Black	7.3%	24.4%	14.5%
	Hispanic	17.1%	2.2%	14.6%
	Others	24.4%	4.4%	12.5%
Education	High school or lower	21.9%	4.5%	4.2%
	Some college	26.8%	28.9%	43.75%
	Bachelor Degree	36.5%	46.6%	31.3%
	Graduate school or higher	14.6%	20%	20.8%
Average Trust with AI systems		2.75	3.26	3.16
Average knowledge level with AI systems		2.53	2.95	2.66

Table A.1: Details of the demographic backgrounds of participants for empirical examinations of the impact of LLM-powered analysis in the income prediction task. *N* represents the number of participants in that treatment.

Demographics		CONTROL (<i>N</i> = 49)	SEQ (<i>N</i> = 40)	ALL (<i>N</i> = 61)
Gender	Male	38.7%	37.5%	31.1%
	Female	57.1%	62.5%	68.9%
	Others	4.2%	0	0
Age	Below 35	47.1%	37.5%	45.9%
	35 -44	38.7%	27.5%	27.4%
	45 or above	14.2%	35%	26.7%
Race	White	73.4%	67.5%	60.7%
	Black	10.2%	12.5%	14.7%
	Hispanic	4.1%	7.5%	6.6%
	Others	12.2%	12.5%	18.0%
Education	High school or lower	16.3%	12.5%	22.3%
	Some college	24.4%	32.5%	23.5%
	Bachelor Degree	36.7%	40%	36.1%
	Graduate school or higher	22.5%	15%	18.1%
Average Trust with AI systems		3.06	2.82	2.88
Average knowledge level with AI systems		2.77	2.72	2.78

Table A.2: Details of the demographic backgrounds of participants for the empirical examinations of the impact of LLM-powered analysis in the recidivism prediction task. *N* represents the number of participants in that treatment.

	Human Solo	AI Solo	Human + AI	GLAD	CATD	LFC	EM	MV
Income Prediction	0.61	0.73	0.76	0.62	0.74	0.61	0.60	0.69
Recidivism Prediction	0.54	0.58	0.61	0.59	0.61	0.58	0.59	0.62

Table B.1: Comparing the accuracy of different combination methods in the income prediction and recidivism prediction tasks, respectively. The best result in each row is highlighted in bold.

feature that the LLM focuses on analyzing in this analysis (e.g., a person’s occupation, a defendant’s charge degree); and (3) the LLM’s analysis on the polarity of this feature towards the final decision, which can be positive, neutral, or negative. For example, in the

recidivism prediction task, a positive polarity indicates that the LLM believes that the feature will increase the probability of reoffending, a neutral polarity indicates that it has no effect, and a negative

polarity indicates that it will decrease the probability of reoffending. $\mu_{\theta_{\text{update}}}$ and $\Sigma_{\theta_{\text{update}}}$ are parameterized by one-layer LSTM-like feedforward networks. The \mathbf{h}^{t-1} is iteratively computed starting from the initial values: $\mathbf{h}^0 \sim \mathcal{N}(\mathbf{h}^0; \mu_{\theta_{\text{init}}}(\mathbf{x}, y^m), \Sigma_{\theta_{\text{init}}}(\mathbf{x}, y^m))$. The dimension of the hidden state \mathbf{h}^t is set as 128. For the final decision making, we employ a two-layer multilayer perceptron to model the decision based on the last hidden state.

C.2 Demographic Information of Participants

The full demographic backgrounds of participants in the study of this phase (where we have five treatments, i.e., CONTROL, SEQ, ALL, ALG, and RANK) for the income prediction and recidivism prediction tasks are shown in Table C.1 and Table C.2, respectively.

D LLM-POWERED ANALYSIS EXAMPLES

D.1 Income Prediction Task

Table D.1 provides additional examples of the GPT-4’s analysis for the income prediction task.

D.2 Recidivism Prediction Task

Table D.2 provides additional examples of the GPT-4’s analysis for the recidivism prediction task.

Demographics		Control (<i>N</i> = 66)	Human-Solo (<i>N</i> = 52)	Seq (<i>N</i> = 85)	All (<i>N</i> = 91)	Rank (<i>N</i> = 53)	Alg (<i>N</i> = 100)
Gender	Male	40.9%	36.5%	41.2%	48.4%	51.8%	52.0%
	Female	57.6%	59.6%	55.3%	51.6%	46.3%	47.0%
	Others	1.5%	3.9%	3.5%	0	1.9%	1.0%
Age	Below 35	49.9%	48.1%	46.9%	30.8%	51.4%	47.0%
	35 -44	31.8%	25.0%	27.1%	33.0%	24.1%	26.0%
	45 or above	18.3%	26.9%	26.0%	36.2%	24.5%	27.0%
Race	White	69.7%	57.6%	61.2%	57.1%	57.4%	67.0%
	Black	10.6%	25%	20%	22.0%	20.3%	15.0%
	Hispanic	9.1%	13.5%	2.4%	9.9%	11.2%	6.0%
	Others	10.6%	3.9%	16.4%	11.0%	11.1%	12.0%
Education	High school or lower	6.1%	3.8%	4.7%	8.8%	9.3%	20.0%
	Some college	30.3%	32.7%	29.4%	28.6%	24.1%	29.0%
	Bachelor Degree	45.5%	30.7%	47.1%	42.9%	51.8%	37.0%
	Graduate school or higher	18.1%	32.8	18.8%	19.7%	14.8%	14.0%
Average Trust with AI systems		2.86	3.55	3.29	3.32	3.29	3.09
Average knowledge level with AI systems		2.43	3.38	2.85	3.03	3.01	2.71

Table C.1: Details of the demographic backgrounds of participants for the evaluation of the effectiveness of the proposed algorithmic selection of LLM-powered analysis in the income prediction task. *N* represents the number of participants in that treatment.

Demographics		Control (<i>N</i> = 84)	Human-Solo (<i>N</i> = 49)	Seq (<i>N</i> = 59)	All (<i>N</i> = 68)	Rank (<i>N</i> = 49)	Alg (<i>N</i> = 88)
Gender	Male	41.6%	43.1%	47.3%	38.2%	15.2%	38.6%
	Female	57.1%	56.9%	50.9%	58.9%	82.6%	60.2%
	Others	1.3%	0%	1.8%	2.9%	2.2%	1.2%
Age	Below 35	41.6%	35.2%	43.8%	44.1%	47.7%	61.3%
	35 -44	34.5%	39.2%	22.80%	33.8%	19.5%	20.4%
	45 or above	23.9%	25.6%	33.4%	22.1%	32.8%	18.3%
Race	White	60.7%	56.8%	71.9%	60.3%	63.1%	62.5%
	Black	16.7%	29.4%	15.8%	14.7%	26.1%	14.7%
	Hispanic	8.3%	7.8%	0	4.4%	6.5%	9.1%
	Others	14.3%	6.0%	12.3%	20.6%	4.3%	13.7%
Education	High school or lower	17.8%	3.9%	10.5%	16.1%	17.4%	5.7%
	Some college	28.6%	29.4%	19.2%	29.4%	34.7%	34.1%
	Bachelor Degree	31.0%	43.1%	49.1%	33.8%	34.8%	46.6%
	Graduate school or higher	22.6%	23.5%	21.2%	20.7%	13.1%	13.6%
Average Trust with AI systems		3.15	3.35	3.00	3.14	2.73	3.03
Average knowledge level with AI systems		2.90	3.47	2.75	2.94	2.63	2.95

Table C.2: Details of the demographic backgrounds of participants for the evaluation of the effectiveness of the proposed algorithmic selection of LLM-powered analysis in the recidivism prediction task. *N* represents the number of participants in that treatment.

Income Prediction	
Gender	<p><u>Value:</u> Female</p> <p><u>Concept:</u> Gender can influence earning potential due to societal and systemic inequalities.</p> <p><u>In this case:</u> As a female, there might be gender pay gaps present, which decreases the likelihood of making over \$50000 per year.</p>
Age	<p><u>Value:</u> 56</p> <p><u>Concept:</u> Middle-aged individuals often have more work experience, which can lead to higher income.</p> <p><u>In this case:</u> At 56, the person likely has significant experience, which increases the likelihood of making over \$50000 per year.</p>
Education Level	<p><u>Value:</u> 10 years</p> <p><u>Concept:</u> Higher education is commonly linked to higher earning potential.</p> <p><u>In this case:</u> With 10 years of education, this might be slightly below the threshold for high-earning positions, which decreases the likelihood of making over \$50000 per year.</p>
Marital Status	<p><u>Value:</u> Divorced</p> <p><u>Concept:</u> Married individuals might have dual incomes or more stability, affecting income levels.</p> <p><u>In this case:</u> Being divorced may lead to financial strain and does not provide the potential dual income advantage of married couples, which decreases the likelihood of making over \$50000 per year.</p>
Occupation	<p><u>Value:</u> Sales</p> <p><u>Concept:</u> Certain professions have higher average incomes than others.</p> <p><u>In this case:</u> Sales positions vary greatly in income, but without further specifics, it's challenging to determine, which has no impact on the likelihood of making over \$50000 per year.</p>
Work Type	<p><u>Value:</u> Self-employed</p> <p><u>Concept:</u> The class of work can affect income level, with self-employed individuals having a high variance in earnings.</p> <p><u>In this case:</u> Being self-employed in sales might offer high earning potential if the business is successful, which has no impact on the likelihood of making over \$50000 per year due to the variability in success amongst self-employed individuals.</p>
Working Hour Per Week	<p><u>Value:</u> 48</p> <p><u>Concept:</u> Working more hours can lead to higher income, especially in jobs with overtime pay or higher rates for longer hours.</p> <p><u>In this case:</u> As a female, there might be gender pay gaps present, which decreases the likelihood of making over \$50000 per year.</p>

Table D.1: Analysis of the Income Prediction task.

Recidivism Prediction	
Gender	<p><u>Value:</u> Female</p> <p><u>Concept:</u> Gender may impact recidivism predictions due to societal perceptions and treatment of different genders within the criminal justice system.</p> <p><u>In this case:</u> Being female, which generally decreases the likelihood of recidivating compared to males, reflecting broader trends in criminal behavior and recidivism.</p>
Age	<p><u>Value:</u> 23</p> <p><u>Concept:</u> Younger individuals are statistically more likely to reoffend, possibly due to lack of stability, economic opportunities, or maturity.</p> <p><u>In this case:</u> At the age of 23, which increases the likelihood of recidivating, as younger age is often associated with higher recidivism rates.</p>
Race	<p><u>Value:</u> Black (African American)</p> <p><u>Concept:</u> Societal and systemic biases related to race can influence recidivism predictions, with some races potentially facing harsher predictions due to historical and ongoing discrimination.</p> <p><u>In this case:</u> Being Black (African American), which might increase the likelihood of recidivating due to systemic biases that affect judiciary outcomes.</p>
Prior Non-juvenile Crimes	<p><u>Value:</u> 0</p> <p><u>Concept:</u> The number of prior criminal charges is a strong predictor of recidivism, with more priors indicating a higher risk.</p> <p><u>In this case:</u> With no non-juvenile criminal charges, which decreases the likelihood of recidivating, suggesting a lack of previous engagement with the criminal justice system.</p>
Juvenile Misdemeanor Crimes	<p><u>Value:</u> 0</p> <p><u>Concept:</u> Juvenile misdemeanors indicate lesser criminal involvement than felonies but can still reflect patterns of behavior leading to recidivism.</p> <p><u>In this case:</u> Having no juvenile misdemeanor charges, which decreases the likelihood of recidivating, representing a lower early involvement in criminal activities.</p>
Juvenile Felony Crimes	<p><u>Value:</u> 0</p> <p><u>Concept:</u> Juvenile felony charges are considered indicators of early criminal behavior, which can predict future recidivism.</p> <p><u>In this case:</u> Having no juvenile felony charges, which decreases the likelihood of recidivating since early criminal behavior is not present.</p>
Charge Issue	<p><u>Value:</u> Driving While License Revoked</p> <p><u>Concept:</u> The specific nature of the current charge can influence recidivism predictions, with certain offenses considered more likely to lead to reoffending.</p> <p><u>In this case:</u> Charged with Driving While License Revoked, which has no impact on the likelihood of recidivating as it may not directly indicate a higher risk of violent or more serious criminal behavior.</p>
Charge Degree	<p><u>Value:</u> Felony</p> <p><u>Concept:</u> The severity of the charge can predict recidivism, with felonies often leading to harsher predictions than misdemeanors.</p> <p><u>In this case:</u> Facing a felony charge, which increases the likelihood of recidivating because felonies are associated with more severe criminal behavior.</p>

Table D.2: Analysis of the Recidivism Prediction task.

REFERENCES

- [1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [2] Gavin Kerrigan, Padhraic Smyth, and Mark Steyvers. 2021. Combining human predictions with model probabilities via confusion matrices and calibration. *Advances in Neural Information Processing Systems* 34 (2021), 4421–4434.
- [3] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [4] Qi Li, Yaliang Li, Jing Gao, Lu Su, Bo Zhao, Murat Demirbas, Wei Fan, and Jiawei Han. 2014. A confidence-aware approach for truth discovery on long-tail data. *Proceedings of the VLDB Endowment* 8, 4 (2014), 425–436.
- [5] Vikas C Raykar, Shipeng Yu, Linda H Zhao, Gerardo Hermosillo Valadez, Charles Florin, Luca Bogoni, and Linda Moy. 2010. Learning from crowds. *Journal of machine learning research* 11, 4 (2010).
- [6] Mark Steyvers, Heliodoro Tejeda, Gavin Kerrigan, and Padhraic Smyth. 2022. Bayesian modeling of human–AI complementarity. *Proceedings of the National Academy of Sciences* 119, 11 (2022), e2111547119.
- [7] Jacob Whitehill, Ting-fan Wu, Jacob Bergsma, Javier Movellan, and Paul Ruvolo. 2009. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise. *Advances in neural information processing systems* 22 (2009).
- [8] Yudian Zheng, Guoliang Li, Yuanbing Li, Caihua Shan, and Reynold Cheng. 2017. Truth inference in crowdsourcing: Is the problem solved? *Proceedings of the VLDB Endowment* 10, 5 (2017), 541–552.