

Coding Negotiations with AI: Instructions and Validation for Coding Model 2

Ray Friedman^{a1}, Jeanne Brett^{b13}, Jaewoo Cho^{a21}, Xuhui Zhan^{a21}, Ningyu Han^{a2},
Sriram Kannan^{a2}, Yingxiang Ma^{a2}, Jesse Spencer-Smith^{a2}, Elisabeth Jackel^{c3},
Alfred Zerres^{c3}, Madison Hooper^{a56}, Katie Babbit^{a6}, Manish Acharya^{a7}, Wendi
Adairⁱⁱ⁴, Soroush Aslani^{g2}, Tayfun Aykac^{pp4}, Chris Bauman^{d4}, Rebecca Bennett^{ff4},
Garrett Brady^{x4}, Peggy Briggs^{o4}, Cheryl Dowie^{e4}, Chase Eck⁴, Igmarr Geiger^{f4},
Frank Jacob^{pp4}, Molly Kern^{cc4}, Sujin Lee^{bb4}, Leigh Anne Liu^{u4}, Wu Liu^{kk4}, Jeffrey
Loewenstein^{v4}, Anne Lytle^{p4}, Li Ma^{w4}, Michel Mann^{s4}, Alexandra Mislin^{mn4}, Tyree
Mitchellⁿ⁴, Hannah Martensen née Nagler^{pp4}, Amit Nandkeolyar^{t4}, Mara
Olekalns^{z4}, Elena Paliakova^{h4}, Jennifer Parlamis^{ij4}, Jason Pierceⁱ⁴, Nancy Pierce^{gg4},
Robin Pinkley^{mm4}, Nathalie Prime^{pp4}, Jimena Ramirez-Marin^{h4}, Kevin
Rockmann^{q4}, William Ross^{o4}, Zhaleh Semnani-Azad^{ee4}, Juliana Schroeder^{r4}, Philip
Smith^{z4}, Elena Stimmer^{pp4}, Roderick Swaab^{j4}, Leigh Thompson^{hh4}, Zhaleh
Thompson^{hh4}, Cathy Tinsley^{k4}, Ece Tuncel^{l4}, Laurie Weingart^{y4}, Robert Wilken^{m4},
JingJing Yao^{dd4}, and Zhi-Xue Zhang^{w4}

¹Project leads.

²Data scientists who built the AI model.

³Negotiation coding experts.

⁴Data contributors.

⁵Statistical analyst.

⁶Coders.

⁷Website designer.

^aVanderbilt University

^bNegotiation and Team Resources (NTR)

^cUniversity of Amsterdam

^dUniversity of California, Irvine

^eUniversity of Houston

^fAalen University

^gUniversity of Wisconsin, Whitewater

^hIESEG School of Management

ⁱUNC, Greensboro

^jINSEAD

^kGeorgetown University

^lWebster University
^mESPC Business School
ⁿLouisiana State University
^oUniversity of Wisconsin, LaCross
^pAnne Lytle & Associates
^qGeorge Mason University
^rUniversity of California, Berkley
^sLeuphana University
^tIndian Institute of Management Ahmedabad
^uGeorgia State University
^vUniversity of Illinois
^wGuanghua School of Management Peking University
^xBocconi University
^yCarnegie Mellon
^zUniversity of Melbourne
^{bb}KAIST
^{cc}City University of New York
^{dd}IESEG School of Management
^{ee}California State University, Northridge
^{ff}University of Central Florida
^{gg}University of North Carolina at Greensboro
^{hh}Norhwestern University
ⁱⁱUniversity of Waterloo
^{jj}University of San Fancisco
^{kk}Hong Kong Polytechnic University
^{mm}Southern Methodist University
ⁿⁿAmerican University
^{pp}ESCP Business School

January 13, 2026

Abstract

This paper provides guidance for scholars wanting to use AI negotiation transcript coding model 2 of the Vanderbilt AI Negotiation Lab. This includes a link to the Lab's website, a description of model 2's codes and coding process, and a report of tests that support the model's validity. It also includes instructions for how to set up and submit your data. When reporting your results from this model, please cite this paper.

Introduction

The Vanderbilt AI Negotiation Lab was created to provide AI tools for negotiation researchers. This includes AI-assisted coding of negotiation transcripts. You can submit your transcripts for coding here: www.AINegotiationLab-Vanderbilt.com. The site provides access to several coding models. For each model, the site provides: definitions of the codes, examples of each code, details of how the AI model was developed and validated, and instructions for how to format and submit your transcripts for coding. That information is also contained in this document.

These models were trained in most cases by learning from transcripts that have been human coded for specific prior projects (a full review of the development of this coding process can be found in Friedman et al, 2024[6]). Each project used a different coding scheme, and even when codes were the same, they may have been interpreted slightly differently. The example sentences provided on the website allow you to see how each project used their codes, so you can decide which model fits your research needs.

We plan to add more coding models and to update models as needed. If you have a set of coded transcripts for a coding scheme that you would like us to include, please contact us.

Model 2 – Brett and Nandkeolyar

Model 2 uses the coding scheme of Brett and Nandkeolyar (unpublished). The study was conducted in 2010, with MBA students from India. The simulation they used was Cartoon (negotiationandteamresources.com). The authors shared with us 43 transcripts they had collected and coded using human coders.

Each of the codes is shown on our website (<https://www.AINegotiationLab-Vanderbilt.com/>). For each of the 12 codes, we provide a definition of the code, a short explanation, and sample sentences from the transcripts. The sample sentences let you see how these scholars operationalized their codes, which is what Model 2 learned from and tries to reproduce when coding your transcripts. As with any coding scheme, different scholars earned from their transcripts and tries might operationalize concepts slightly differently. You should decide if this coding scheme will be useful to you by reviewing how the authors used it.

Coding Process

Transcripts are coded in three steps.

1. **Unitization (you need to do this).** The model provides one code for each set of words or sentences that you identify as a unit in your Excel document. You can choose to have units be speaking turn, sentences, or thought units. The easiest to set up is speaking turns, since switching between speakers is clearly identifiable in transcripts. The next easiest is sentences, since they are identified by one of a period, question mark, or exclamation mark. However, different transcribers may end sentences in different places. The hardest unit to create is the thought unit since that takes careful analysis and can represent as much work as the coding itself. (See the NegotiAct coding manual for how to create thought units; Jackel et al, 2022[10]). Clarity of meaning runs the opposite direction. The longer the unit, the more likely there are multiple ideas in the unit, and less clarity for human or AI coders to know what part to code. Brett and Nandkeolyar coded speaking turns, but 68% of their speaking

turns contained just one sentence. The closest alignment with the training data would be for you to use sentences as the unit.

2. **Model Assigns Codes.** The model assigns a code to each unit you submit, based on in-context learning. Coding is guided by the prompt we developed and tested. For more on in-context learning see Xie and Min (2022)[14]. Our prompt for this model includes several elements:
 - (a) Five fully coded transcripts. These transcripts were chosen from the 43 available transcripts in the following way. First, any combination of five was considered only if that set included all 12 codes. Second, five of those combinations were chosen at random to test. Third, the one that produced the highest level of match with human coders was retained.
 - (b) Instructions to pay attention to who was speaking, such as “buyer” or “seller”.
 - (c) Instructions to pay attention to what was said in the conversation before and after the unit being coded.
 - (d) Supplementary instructions about the difference between “substantiation” and “information” since in early tests the model often coded substantiation as information, and vice versa. This confusion is not surprising since substantiation usually comes in the form of providing information, but with the purpose of supporting a specific offer or demand.
 - (e) Additional examples of any codes where the five training transcripts did not contain at least 15 examples. We created enough additional examples (based on our understanding of the code) to bring the examples up to 15. We needed to add 12 examples of multi-issue offer, 13 examples of offer accepted, 12 examples of process comments and 14 examples of Miscellaneous Off-Task.
3. **We Run the Model Five Times.** We automatically run the model five times, to assess consistency of results. As expected the results are not always the same, since with in-context learning the model learns anew with each run and may learn slightly differently each time. Variation is also expected since some units may reasonably be coded in several ways. By running the coding model five times, we get five codes assigned to each speaking unit. If three, four, or five of the five of the runs have the same code, we report the code and indicate the level of “consistency” of that code (three, four, or five out of five). If there are not at least three consistent results out of five runs, or if the model fails to assign a code, we do not report a model code. In these cases, the researcher needs to do human coding.

Validation of Model

This section reports the initial validation, using Claude Sonnet 3.5. This version of Sonnet was decommissioned fall, 2025, and replaced in our model with Sonnet 4.5. Updated validation of key measures is reported at the end of this section and shows equivalent results for Sonnet 4.5.

Validation occurred in several steps.

Step 1: Compare the model coding with humans by Brett and Nandkeolyar. To do this, we asked the model to code the 3496 units contained in the Brett and Nandkeolyar transcripts that were not selected for training. We looked at several criteria.

- (a) **Consistency.** In our test, there was “perfect” consistency of model coding (five out of five runs of the model assigned the same code to a unit) for 3,448 of the units, “high” consistency (four out of five) for 48 of the units, “modest” consistency (three out of five) for zero of the units, and no cases where the model did not report a code. Thus, 98.6% of codes had perfect consistency.
- (b) **Match with human coding.** We assessed whether the model assigned the same code as the human coders. The overall match level for units where the model assigned a code was 75% (95% CI: .74, .76). To ensure that the model was not biased by matching more accurately with the human coder in early or later phases of the negotiation, we tested whether the match level was different for coding the first versus second half of all transcripts. The match level was 75% for the first half and 74% for the second half, suggesting no bias based on phase of the negotiation. We also looked at match by level of consistency (see Table 2). These results suggest that users may want to accept model-assigned codes only for those codes where the model achieves perfect consistency (five out of five).

Table 1: Match Percentage by Consistency Level, Validation Step 1

Level of Consistency	Match with Human Codes	% Achieve This Consistency Level Among Those Assigned a Code	Number	Match	Percentage Match
Modest Consistency	3 out of 5	0%	0	not match	
			0	match	n/a
High Consistency	4 out of 5	1.4%	29	not match	
			19	match	40%
Perfect Consistency	5 out of 5	98.6%	845	not match	
			2,603	match	75%

*0 cases did not reach the 3 out of 5 consistency threshold or the model failed to assign a code

We also calculated the Cohen’s kappa, with the model codes as coming from one rater and the human coding as coming from a second rater. This calculation, compared to the “percentage match,” accounts for matches that might occur based on chance. Cohen’s kappa was calculated in R (R Core Team, 2022) using the IRR package (Gamer & Lemon, 2019[7]). Cohen’s kappa was equal to 0.70, with the no information rate of .26 (p-value of difference is < .001). According to Landis and Koch (1977) this represents “substantial agreement”, and according to Fleiss (1981)[5] is “fair to good” agreement. Rather than relying on conventional categorical guidelines to interpret the magnitude of kappa, Bakeman (2023) argues that researchers should estimate observer accuracy or how accurate simulated observers need to be to produce a given value of kappa. The KappaAcc program (Bakeman, 2022[4]) was used to estimate observer accuracy, which was found to be 86%.

It is also worth noting that in many cases where scholars establish inter-coder reliability, there is a process of cross-rater discussion that is used to resolve initial

differences of opinion between the two coders. In a study of inter-coder agreement, coder agreements in the 80% range began with initial coder agreements in the 40% range (Garrison, et al 2005[8]). Of course, in our case with initial coder there can be no cross-rater discussion between a model and a human, taking away one step that is often used to achieve higher kappas. The closest we can get to that process is to have a third person view the cases of human-model disagreement to provide a judgment of which code was more correct. Also, the fact that so many codes need human-to-human discussions to resolve, suggests some inherent ambiguity about code assignments and opens up the possibility that several different codes might reasonably be assigned to some segments of transcripts.

- (c) **Summary Data and Confusion Matrix.** We created a confusion matrix for all codes (see see Figure 1). The vertical axis shows human coding. The horizontal axis shows model coding. Also included below (see see Table 2) are summary statistics showing which codes appeared most frequently in the human coding (Substantiation was most common representing 31.64% of the codes, while Miscellaneous Off-Topic was least common representing just .11% of the codes), and the human-model match level for each code. The highest levels of human-model match (other than for Miscellaneous) were for Positive Response, Question, and Multi-Issue Offer. There appears to be a rough correlation between number of units and match percentage, suggesting that match percentage goes up when there are more examples of that code in the training transcripts for the model to learn from and when there are more opportunities to find that code in the test transcripts.

In terms of absolute numbers of mismatches, the largest set is 155 human-coded Substantiation codes that were coded as Information by the model. This is an issue we recognized early in our testing, which resulted in added instructions in the prompt to reduce this mismatch. The fundamental problem is that Substantiation is often achieved by providing information, but to be Substantiation that information must support a particular argument or claim. There were also 43 cases of human-coded Miscellaneous On-Task that were coded as Information by the model. The next largest set of mismatches were 42 where humans assigned a code of Single issue offer while the model assigned a code of Multi-issue offer, which is easy to imagine happening.

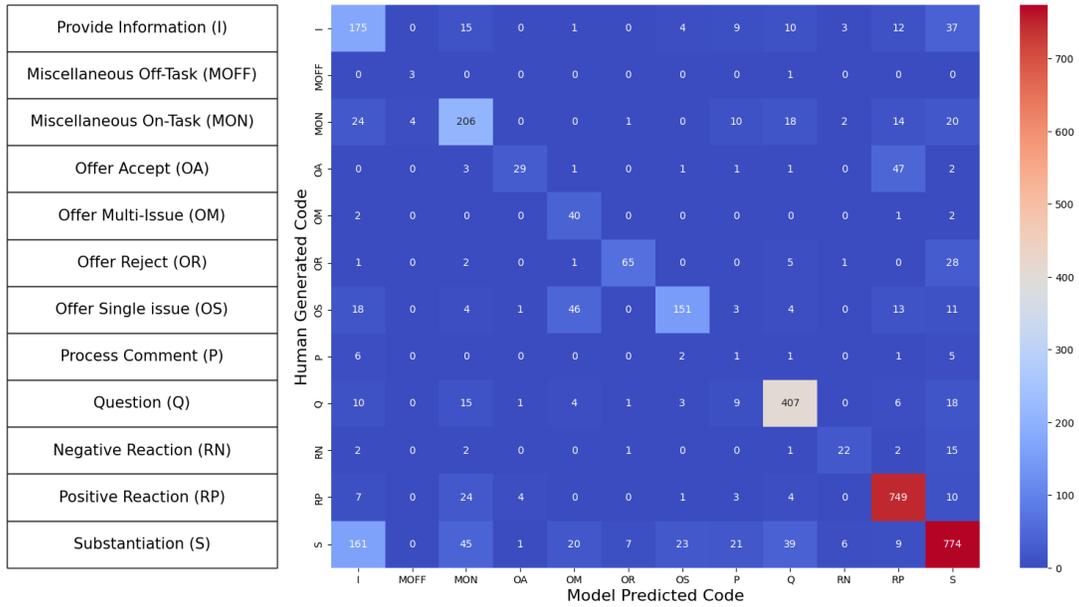


Figure 1: Confusion Matrix, Validation Step 1

Human Code	% of units Across All Transcripts	Model Match %
S	31.64	70%
RP	22.94	93%
Q	13.56	86%
MON	8.55	69%
I	7.61	66%
OS	7.18	60%
OR	2.95	63%
OA	2.43	34%
RN	1.29	49%
OM	1.29	89%
P	0.46	6%
MOFF	0.11	75%

Table 2: Mach Percentage by Code, Validation Step 1

Closer Look at Mismatches. To assess the nature of these and other mismatches, we selected a random sample of 100 mismatches for closer examination. Given that original human coders may be just as likely to make errors (or simply vary in their judgments) as the model, we wanted to see if newly trained coders would see the human codes or the model-provided codes as more accurate. We provided

these coders with the 100 speaking turns, as well as the two speaking turns preceding that speaking turn, along with the human and model codes. They were not informed which code came from the model or humans, and the order in which they saw the two codes was flipped halfway through the 100 samples to avoid order effects. They selected which code they saw as more accurate. This was done first separately by the two coders, and then they were asked to resolve through discussion any cases where they disagreed. In the end, these new coders thought the model-provided codes were correct 67% of the time and the human-generated codes 25% of the time, and were uncertain about which was correct 8% of the time. Based on this we can expect that the model is correct in 67% of the cases with mismatches, so we can trust that about 92% of the model codes are accurate.

Step 2: Match with Human coding for Different Simulations.

The first step of validation involved matching human and model codes where the negotiation simulation used for training was the same as the negotiation simulation used for testing the model (Cartoon). But users may have transcripts from any number of simulations or real-world negotiations, not just the simulation used in Brett and Nandkeolyar (unpublished) study. Therefore, we wanted to test how well the model would match human coders who applied the Brett and Nandkeolyar model to transcripts using other simulations. We selected a set of 3 Transcripts from a study that used The Sweet Shop simulation, and 3 transcripts from a study that used the Les Florets simulation. Since these transcripts were not initially coded using the Brett and Nandkeolyar codes, we needed to train two coders to use the Brett and Nandkeolyar codes. After initial training, they reached a level of inter-coder reliability of $k=.73$. They coded the transcripts separately and came together to discuss any cases where they disagreed and assign a code. This provided the human codes for a set of the Sweet Shop and LesFlorets simulations. These transcripts were then coded using our model.

The 6 transcripts had 1302 speaking turns, of which 99% were single sentences. The model had perfect consistency for 87% of the speaking turns (all five runs assigned the same code), high consistency for 10% of the speaking turns (4 out of 5 runs assigned the same code), and modest consistency for 1.6% of the speaking turns (3 out of 5 runs assigned the same code). There was 1 case of less than 3 out of 5 consistency. The match percentage was 67% for perfect consistency codes, 81% for high consistency codes, and 45% for moderate consistency codes (see Table 3). Overall, the match percentage was 68%. This was lower than our prior tests, as expected, because these transcripts did not have the same issues and topics as training transcripts (which used The Cartoon simulation). For that reason, these results may better represent the model’s effectiveness with most transcripts. We also checked to see if one set of transcripts did better than the other. The match percentage was also 70% for just the Les Florets transcripts and 64% for just the Sweet Shop transcripts, suggesting that the model should do just as well with transcripts using different simulations.

We also calculated the Cohen’s kappa. The weighted Cohen’s kappa was .63 with the no information rate of .26 (p-value of difference is $< .001$). This kappa according to Landis and Koch (1977)[11] is “moderate agreement”, and according to Fleiss (1981) is “fair to good” agreement. Rather than relying on conventional categorical guidelines to interpret the magnitude of kappa, Bakeman (2023)[4] argues that researchers should

estimate observer accuracy or how accurate simulated observers need to be to produce a given value of kappa. The KappaAcc program (Bakeman, 2022)[4] was used to estimate observer accuracy, which was found to be 81%.

Table 3: Match Percentage by Consistency Level, Validation Step 1

Level of Consistency	Match with Human Codes	% Achieve This Consistency Level Among Those Assigned a Code	Number	Match	Percentage Match
Modest Consistency	3 out of 5	.02%	12	not match	45%
			10	match	
High Consistency	4 out of 5	10.1%	25	not match	81%
			106	match	
Perfect Consistency	5 out of 5	87.8%	377	not match	67%
			771	match	

*1 case did not reach the 3 out of 5 consistency threshold or the model failed to assign a code

The proportion of speaking units that fell into each category were roughly similar to what we saw in the first validation tests, with most speaking units being: Information, Question, and Response Positive. In this set of transcripts Substantiation was also fairly common (see Table 2). As with the first validation test, model-human match percentage appears to be highly correlated with number of codes.

The confusion matrix (see Figure 2) shows that, once again, the largest number of mismatches comes from Information/Substantiation. It also shows that nearly all of the mismatches were cases where the model assigned a code of “information” when the humans assigned various other codes.

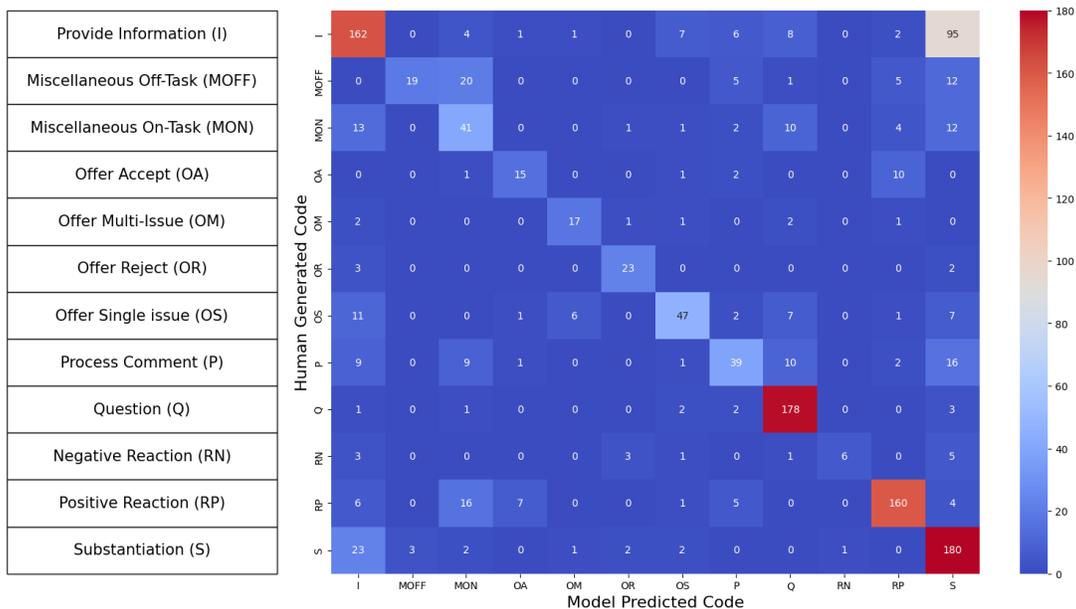


Figure 2: Confusion Matrix, Validation Step 1

Human Code	% of units Across All Transcripts	Model Match %
I	21.98	57%
S	16.45	84%
RP	15.30	80%
Q	14.37	95%
P	6.69	45%
MON	6.46	49%
OS	6.30	57%
MOFF	4.77	31%
OA	2.23	52%
OR	2.15	82%
OM	1.84	71%
RN	1.46	32%

Table 4: Mach Percentage by Code, Validation Step 1

In order to assess the mismatches, we collected a random sample of 98 sentences with mismatches, along with the two prior sentences and the human and model codes. We then took off the column labels and randomly mixed the order of the codes. Since the human coding in this case was done by our coding team, we wanted a different person to select which of the two codes was more correct. This was done by the first author. The

results are shown in Table 5. About 47% of the time the human code was deemed more accurate, while in 40% of the cases the model was deemed more accurate. In another 6% of the cases, both were deemed correct (because, for example, the sentence was long and really contained two thought units). In 7% of the cases both the human and model codes were deemed incorrect or the sentence was deemed uninterpretable (because, for example, there were missing words in the transcription). Looking, then, at the 32% of speaking units that were a mismatch, perhaps half of them might still be deemed accurate, bringing the match percentage up from 68% to about 85%.

Code Selection		Count	
Clear Choice	Human Code is Correct	46	85
	Model Code is Correct	39	
Both Correct		6	6
Both Incorrect	Human and Model Both Incorrect	2	2
Not Understood	Could not Understand the Sentence	5	5

Table 5: Assessment of 98 Sample Mismatches

Confirm Validation of Model Using Clause Sonnet 4.5 (replaced Sonnett 3.5 in fall, 2025)

This section reports the initial validation, using Claude Sonnet 3.5. This version of Sonnet was decommissioned fall, 2025, and replaced in our model with Sonnet 4.5. Here we report re-runs of the main analyses which show that the Sonnet 4.5 validation results are consistent with the validation results reported above for Sonnet 3.5.

Step 1: Confirm Validation. Compare the model coding (using Sonnet 4.5) with humans by Brett and Nandkeolyar (unpublished). We asked the model to code the first half of all the transcripts in Brett and Nandkeolyar (unpublished) that were not selected for training, or 1526 speaking units. The overall match level for units where the model assigned a code was 73%, which is very close to the match results (75%) found when we used Sonnet 3.5. We calculated Cohen’s kappa, with the model as one rater and the human coders as a second rater, using the KappaAcc program (Bakeman, 2022). Cohen’s Omnibus Kappa was .67 (nearly identical to the .69 Kappa when using Sonnet 3.5), with the no information rate of .26 (p-value of the difference is $\leq .001$). Estimated observer accuracy was 84% (nearly identical to the 85% estimated when using Sonnet 3.5).

Step 2: Confirm Validation. Compare the model coding (using Sonnet 4.5) with human coding from Brett and Nandkeolyar (unpublished). We used the same 6 transcripts used in Step 2 analysis for Sonnet 3.5 reported above, which contained 1301 speaking turns. The match percentage was .65 using Sonnet 4.5, close to the .68 results using Sonnet 3.5. We calculated Cohen’s kappa, with the model codes as one rater and the human codes as a second rater, using the KappaAcc program (Bakeman, 2022). Cohen’s Omnibus Kappa was .60 (close to the .63 Kappa when using Sonnet 3.5), with the no information rate of .23 (p-value of the difference is $\leq .001$). Estimated observer accuracy was 80% (nearly

identical to the 79% estimated when using Sonnet 3.5). These results show that the coding performance of Sonnet 4.5 is nearly identical to that of Sonnet 3.5.

Formatting Your Transcripts

Set up your transcripts for analysis by putting them into an excel sheet. The format should be as shown below. Label the first column “SpeakerName” and list whatever names you have for those speakers (e.g., buyer/seller, John/Mary). Label the second column “Content” and include the material (in English or other major languages) that is contained in your unit of analysis (which may be a speaking turn, a sentence, or a thought unit). Also include columns for ResearcherName, Email, and Institution and include that information in the next row. If you use speaking turns then speakers will alternate, and the format will look like this:

SpeakerName	Content	ResearcherName	Email	Institution
<i>Buyer</i>	<i>All words in speaking turn...</i>	<i>Your name</i>	<i>Your email</i>	<i>Your institution</i>
<i>Seller</i>	<i>All words in speaking turn...</i>			
<i>Buyer</i>	<i>All words in speaking turn...</i>			
<i>Seller</i>	<i>All words in speaking turn...</i>			
<i>Buyer</i>	<i>All words in speaking turn...</i>			
etc	etc			

If you use *sentences* or *thought units* then it is possible that speakers may appear several times in a row, and the format will look like this:

SpeakerName	Content	ResearcherName	Email	Institution
<i>Buyer</i>	<i>All words in sentence or thought unit...</i>	<i>Your name</i>	<i>Your email</i>	<i>Your institution</i>
<i>Buyer</i>	<i>All words in sentence or thought unit...</i>			
<i>Buyer</i>	<i>All words in sentence or thought unit...</i>			
<i>Seller</i>	<i>All words in sentence or thought unit...</i>			
<i>Seller</i>	<i>All words in sentence or thought unit...</i>			
<i>Buyer</i>	<i>All words in sentence or thought unit...</i>			
etc	etc			

Create one Excel file for each transcript. Name each file in the following way:

YourName_StudyName_1

YourName_StudyName_2

YourName_StudyName_3

Etc.

Submit Your Transcript

To submit your transcript for the model to code, drag and drop one or several transcript files into the submission section of the website. The content can be in English or other major languages. It will likely take about 10 minutes for Claude to process each transcript, although this can vary based on how much demand Claude has at the moment you submit your files. **Do not close your window while you are waiting for results - you will lose your results.** Once the analysis for each transcript is complete, you will receive the output in a csv file that is automatically downloaded to your download folder. We suggest submitting just a few files at a time, so that you can check the output before doing too many analyses. The output files will include:

- Transcript Name
- Speaker
- The text (which could be a thought unit, sentence, or speaking turn)
- The code assigned to that text
- Consistency score for that code

If you have any questions, contact the Vanderbilt AI Negotiation Lab.

References

- [1] Aslani et al. (2014), "Measuring negotiation strategy and predicting outcomes: Self-reports, behavioral codes, and linguistic codes," presented at the annual conference of the International Association for Conflict Management, Leiden, The Netherlands.
- [2] Aslani, S., Ramirez-Marin, J., Brett, J., Yao, J., Semnani-Azad, Z., Zhang, Z. X., ... & Adair, W. (2016). Dignity, face, and honor cultures: A study of negotiation strategy and outcomes in three cultures. *Journal of Organizational Behavior*, 37(8), 1178-1201.
- [3] Adair, W. L., & Brett, J. M. (2005). The negotiation dance: Time, culture, and behavioral sequences in negotiation. *Organization Science*, 16, 33-51.
- [4] Bakeman, R. (2022). KappaAcc: A program for assessing the adequacy of kappa. *Behavior Research Methods*. <https://doi.org/10.3758/s13428-022-01836-1>
- [5] Fleiss, J.L. (1981). *Statistical methods for rates and proportions* (2nd ed.). New York: John Wiley. ISBN 978-0-471-26370-8.
- [6] Friedman, R., Brett, J., Cho, J., Zhan, X. et al. (2024). An application large language models to coding negotiation transcripts: The Vanderbilt AI negotiation lab. (forthcoming)
- [7] Gamer, M., Lemon, J., Fellows, I., & Singh P. (2019) irr: Various coefficients of interrater reliability and agreement. R package version 0.84.1. <https://CRAN.R-project.org/package=irr>
- [8] Garrison, D. Cleveland-Innes, M., Koole, M. & Kappelman, J. (2006). Revisiting methodological issues in transcript analysis: Negotiated coding and reliability. *The Internet and Higher Education*. 9. 1-8. 10.1016/j.iheduc.2005.11.001.

- [9] Gunia, B. C., Brett, J. M., Nandkeolyar, A. K., & Kamdar, D. (2011). Paying a price: Culture, trust, and negotiation consequences. *Journal of Applied Psychology*, *96*, 774-789.
- [10] Jackel, E., Zerres, A., Hamshorn de Sanchez, C., Lehmann-Willenbrock, & N., Huffmeier, J. (2022), "NegotiAct: Introducing a comprehensive coding scheme to capture temporal interaction patterns in negotiations," *Group and Organization Management*. (See supplementary file for coding guidelines.)
- [11] Landis, J.R. & Koch, G.G. (1977). The measurement of observer agreement for categorical data. *Biometrics*. *33* (1): 159-174. doi:10.2307/2529310.
- [12] R Core Team (2022). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>.
- [13] Weingart, L. R., Thompson, L. L., Bazerman, M. H., & Carroll, J. S. (1990). Tactical behavior and negotiation outcomes. *International Journal of Conflict Management*, *1*, 7-31.
- [14] Xie. S.M. & Min, S. (2022). How does in-context learning work? A framework for understanding the differences from traditional supervised learning. Stanford AI Lab Blog, Aug 1. <https://ai.stanford.edu/blog/understanding-incontext/>