
Factorized embedding of goal and uncertainty in the lateral prefrontal cortex guides stably flexible learning

Received: 4 December 2024

Accepted: 12 November 2025

Cite this article as: Sung, Y., Rigotti, M., Lee, S.W. Factorized embedding of goal and uncertainty in the lateral prefrontal cortex guides stably flexible learning. *Nat Commun* (2025). <https://doi.org/10.1038/s41467-025-66677-w>

Yoondo Sung, Mattia Rigotti & Sang Wan Lee

We are providing an unedited version of this manuscript to give early access to its findings. Before final publication, the manuscript will undergo further editing. Please note there may be errors present which affect the content, and all legal disclaimers apply.

If this paper is publishing under a Transparent Peer Review model then Peer Review reports will publish with the final article.

October 22, 2025

Daniel Barry, Ph. D.
Senior Editor,
Nature Communications

Dear Dr. Barry,

We are pleased to submit the final revised version of our manuscript, "Factorized embedding of goal and uncertainty in the lateral prefrontal cortex guides stably flexible learning" (Manuscript ID: NCOMMS-24-79943B-Z), for publication in *Nature Communications*.

As requested, we have carefully addressed all editorial points outlined in the Author Checklist. We have uploaded all the required files with this submission, including the completed checklist, the revised manuscript, separate figure files, single-PDF Supplementary Information, Source Data, and a revised Reporting Summary.

Brief summary:

Behavior and fMRI reveal that the human lateral prefrontal cortex factorizes goal and uncertainty into a geometric code, balancing flexible goal pursuit with stable control to prevent erratic behavior.

Thank you again for considering our final revision.

Yours sincerely,



Sang Wan Lee, Ph. D.

Director, Center for Neuroscience-inspired AI
Associate Professor,
Department of Brain and Cognitive Sciences,
Department of Bio and Brain Engineering,
Kim Jaechul Graduate School of AI,
Graduate School of Data Science,

Korea Advanced Institute of Science and Technology (KAIST)

291 Daehak-ro, Yuseong-gu, Daejeon 34141, Republic of Korea
Email: sangwan@kaist.ac.kr Web: <http://aibrain.kaist.ac.kr>

Author Checklist

Manuscript Number:

NCOMMS-24-79943B-Z

Please check the items below carefully and add a response in each row of the table to indicate the changes that you have made. Please also check through any additional marked-up edits we may have provided within the manuscript file.

Remaining reviewer comments

Our guidance:

Please ensure that the response to Reviewer #1 regarding the conceptual advance over your prior work is clearly articulated in the manuscript itself and not just in the response letter.

Your response:

Our key conceptual advance is the framing of the stability-flexibility trade-off in goal-directed learning and explaining it at the level of neural representation. This core question is established from the introduction and is logically supported by our results and figures. This framework has been strengthened throughout the revision process by incorporating the reviewers' valuable feedback on the task design, as well as the behavioral and neural analyses.

We ensured to fully articulate the conceptual advance in the manuscript, we added a paragraph in the previous-round revision to directly compare our findings with prior work (second paragraph in the Discussion). To further clarify this distinction and proactively address similar questions from future readers, we have updated this paragraph as follows:

"Our study presents a significant conceptual advance by offering a representational solution to the stability-flexibility dilemma in goal-directed learning. While prior works have primarily focused on identifying 'what' variables are encoded in the PFC (e.g., value, uncertainty), our work elucidates 'how' the brain navigates this trade-off through the specific geometrical structure of its neural representations. Earlier studies in decision-making have attempted to probe behavioral flexibility in response to changes in goals or context^{citep{FeherDaSilva.etal2023, Daw.etal2011, Saez2015, Mohring.Glascher2023, Kool.etal2017, Kim.etal2024}}, while placing less emphasis on the concurrent challenge of maintaining stability against environmental noise. Additionally, studies on uncertainty representation have often centered on perceptual judgments^{citep{Kiani_Shadlen_2009, Hebart.etal2016, Gherman_Philastides_2018, vanBergen.etal2015, Li.etal2021, Geurts.etal2022}} rather than on the complex dynamics of sequential action toward delayed goals. Our findings bridge this gap by showing, through the lens of representational geometry, that the PFC employs a joint, factorized coding scheme for goals and uncertainty. This neural architecture offers a mechanistic account of how the brain reconciles competing cognitive demands, enabling robust and generalizable goal pursuit while remaining attuned to environmental statistics, thus supporting both adaptive behavior and resilience to noise."

Author information

Our guidance:

We ask that you consult with your coauthors to ensure that all names, affiliations, and titles are represented correctly. Note that if any authors are added or removed after this point then all authors will be requested to provide approval documentation that could potentially delay the production of your paper.

Your response:

All authors have reviewed and confirmed the accuracy of names, affiliations, and titles.

Ensure affiliations are appropriately labeled and featured sequentially and in ascending order (1,2,3,... or a,b,c,...). Please ensure all corresponding authors are marked with a specific symbol and include their emails. Similarly, if you have "equally contributing" or "joint supervision" authors, use a specific symbol to mark them and not a number.	Affiliations are labeled sequentially and the corresponding author is marked with a symbol, with an email address included.
Please ensure the author contributions section mentions each author's initials at least once with their contributions to the work. Authors with the same initials must be differentiated in the statement.	The Author Contributions statement lists each author's initials at least once.

Article structure

Our guidance:

Your response:

We can accommodate up to 10 display items (Figures or Tables) in the main article. Each Figure and Table must fit easily within an A4 page (210 x 297 mm). Please ensure that the number and size of your Figures and Tables fulfil these requirements to avoid any delay in the acceptance of your article.	We have fewer than ten main figures, and each fits within a single A4 page.
Ensure main Figures are uploaded as separate individual files. Each figure file must contain all intended panels labelled and displayed as intended and fit entirely on a single page. Do NOT include legends within the figure files, as these must be in the main manuscript. Supplementary Figures must be all contained in the Supplementary Information PDF and do NOT need to be uploaded separately.	We confirmed that all the figures fulfilled the requirements.
Please ensure your main manuscript file includes the following sections, in this order: <i>Title Author list Affiliations Abstract Introduction Results Discussion (optional) Results and Discussion (optional) Methods Data Availability Code Availability (if relevant) References Acknowledgements Author Contributions Statement Competing Interests Statement Tables Figure Legends/Captions (for main text figures)</i>	The main manuscript follows the required section order.
We do not edit Supplementary Information files; they will be uploaded with the published article as they are submitted with the final version of your manuscript. Any tracked changes should be removed from the file and the file should be provided as a PDF file. Supplementary Figures do not need to be provided separately. Please supply Source Data files for all data presented in graphs within the Figures.	Supplementary Information is provided as a single clean PDF. Source Data are supplied in an Excel file with one worksheet per figure.
Within the Source Data file, the relevant raw data from each figure or table (in the main manuscript and in the Supplementary Information) should be represented by a single sheet in an Excel document, or a single .txt file or other file type in a zipped folder. An example of the Source Data file is available demonstrating the correct format: https://www.nature.com/documents/ncomms-example-source-data.xlsx The file should be labelled 'Source Data', with the title and a brief description included in your response here, and should be mentioned in all relevant figure legends using the template text below: 'Source data are provided as a Source Data file.' A reference to the source data file should be added in the 'Data Availability' section, using the text: Source data are provided with this paper.	The file is titled "Source Data.xlsx" with one worksheet per figure. The phrases "Source data are provided as a Source Data file." (in relevant figure legends) and "Source data are provided with this paper." (in Data Availability) have been added.

Main text

Our guidance:

All references to frequentist inferential statistics must be reported as statistic(degrees of freedom) = value, p = value, effect size statistic = value, % Confidence Intervals = values

Your response:

We reported the effect size or % confidence intervals appropriately in the main text, and additionally the statistic(df) and p-value where it is required. All the statistical details, statistic values and p-values are reported in the figure legends and Supplementary Tables 1-4.

Please do not use italics, bold font, underlining or speech marks/quotation marks except in headings unless required for technical terms (in both the main text and the display items).

Italics or quotation marks are only used for some technical terms (e.g. *goal condition* and *uncertainty condition*, *choice versatility* and *choice consistency*).

Please make sure that mathematical terms throughout your manuscript and Supplementary Information (including in figures, figure axes, and legends) conform strictly to the following guidelines. Equations must be supplied in editable format, and not as images. Scalar variables (e.g. x, V, x) must be typeset in italic, whereas multi-letter variables and functions (e.g. log) must be formatted in roman. Vectors (such as the wavevector k or the magnetic field vector B) must be typeset in bold without italics.

All requirements were fulfilled.

Please label equations sequentially as (1), (2), (3), etc.

The manuscript does not include in-text citation of equations, so the equations are not labelled.

Figures and Tables

Our guidance:

Please see the guidelines linked below for detailed instructions about how your figures should be prepared. Following these instructions will reduce the chances of delays should we need to request replacement artwork from you at a later stage.

<https://www.nature.com/documents/NRJs-guide-to-preparing-final-artwork.pdf>

Your response:

We have followed the artwork guidelines and provided editable vector files (.ai, .pptx).

Please ensure that data presented in a plot, chart or other visual representation format shows data distribution clearly (e.g. dot plots, box-and-whisker plots, violin plots). When using bar charts, please overlay the corresponding data points (as dot plots) whenever possible and always for $n \leq 10$. All box-plot elements (center line, limits, whiskers, points) should be defined in the legends accompanied by precise n numbers.

Figures 2b, 2d, 3b, 4a have been revised to visualize data distributions to comply the policy (bars with overlaid data points).

Please note that data presentation has to be revised to comply with our policy in figure(s) 2b, 2d, 3b, 4a

All error bars need to be defined in the legends (e.g. SD, SEM) together with a measure of centre (e.g. mean, median). For example, the legends should state something along the lines of "Data are presented as mean values +/- SEM" as appropriate. All box plots need to be defined in the legends in terms of minima, maxima, centre, bounds of box and whiskers and percentile.

We added 'Data are presented as mean \pm SEM from $n=20$ participants.' to the figure 2, 3b, and 4a legends.

Please note that the error bars/error bands need to be defined in the legend(s) of figure(s) 3b, 4a

Please note that the measure of centre for the error bars/error bands needs to be defined in the legend(s) of figure(s) 2a, 2c

The figure legends must indicate the statistical test used. Where appropriate, please indicate in the figure legends whether the statistical tests were one-sided or two-sided and whether adjustments were made for multiple comparisons. For null hypothesis testing, please indicate the test statistic (e.g. F, t, r) with confidence intervals, effect sizes, degrees of freedom and P values noted. Please provide the test results (e.g. P values) as exact values whenever possible and with confidence intervals noted.

Please indicate the statistical test used for data analysis and where appropriate, please specify whether it was one-sided or two-sided and whether adjustments were made for multiple comparisons, in the legend(s) of figure(s) 1f, 3c, 3d, 4b	Corresponding indications were added as below: 1f legend: ... For simulation results, each point represents an agent, and MB and MF agents were simulated with 20,000 sets of random parameters, respectively (4.4). See Supplementary Table 1 for the statistical details of the Pearson's correlation test (two-sided). ... 3c, d legend: ... Only statistically significant correlations are represented with filled bars (Pearson's correlation, *: p < 0.05, **: p < 0.01, ***: p < 0.001; two-sided test). For all the exploratory correlation analyses, multiple comparison corrections were applied with a false-discovery rate (Benjamini–Hochberg procedure) for the number of ROIs with q = 0.05. ... 4b legend: ... Only statistically significant correlations are represented with filled bars (Pearson's correlation, *: p < 0.05, **: p < 0.01, ***: p < 0.001; two-sided test), corrected for the number of ROIs (Benjamini–Hochberg procedure, q=0.05). ...
Tables may not contain colour - please reformat your tables to black and white	No main-text tables are included.
Shadings or symbols in graphs must be defined in some fashion. We prefer that you use a key within the image; do not include colored symbols in the legend/caption.	The requirement is fully fulfilled.
Any abbreviations, symbols or colours present in your figures must be defined in the associated legends.	The requirement is fully fulfilled.
Please ensure that any abbreviations (for example brain regions) used in the figures are defined in the figure legends.	All abbreviations (e.g., brain regions) are defined at first appearance in each figure legend.
Please remove "previous page" from the figure titles.	It was removed from the Fig. 3 title.

Data and Code

Our guidance:

Nature journals strongly support public availability of data and code. Please deposit the data and code used in your paper into a public data repository, or alternatively, present the data as Supplementary Information. If data can only be shared on request, please explain why in your Data Availability Statement, and also in the correspondence with your editor.

Please note that for some data types, deposition in a public repository is mandatory. Any restrictions on sharing of these data types must be clearly indicated in the statement and discussed with the editor. More information on our data deposition policies and available repositories can be found here:

<https://www.nature.com/nature-research/editorial-policies/reporting-standards#availability-of-data>

Your response:

Data and code are publicly available: OSF (<https://osf.io/2gyue>) and GitHub (<https://github.com/brain-machine-intelligence/RLdim-mvpa-model>), as referenced in the Data/Code Availability sections.

<p>All published manuscripts reporting original research in Nature Portfolio journals must include a data availability statement, within the Methods and under the heading 'Data Availability'.</p> <p>The data availability statement must make the conditions of access to the "minimum dataset" that are necessary to interpret, verify and extend the research in the article, transparent to readers. We ask that you don't use phrases like 'available on reasonable request' but instead specify any restrictions to accessing your data as described below.</p> <p>This minimum dataset may be provided through deposition in public community/discipline-specific repositories, custom proprietary repositories or general repositories like Figshare, Zenodo and Dryad. Providing large datasets in supplementary information is strongly discouraged and the preferred approach is to make data available in repositories. Please see https://www.springernature.com/gp/authors/research-data-policy/recommended-repositories for a list of recommended repositories.</p> <p>If DOIs are provided, we also strongly encourage including these in the Reference list (authors, title, publisher (repository name), identifier, year).</p> <p>The Data Availability Statement should also reference any source data published alongside the paper.</p> <p>For clinical datasets or third party data, please ensure that the Data Availability statement adheres to our policy (https://www.nature.com/nature-research/editorial-policies/reporting-standards#availability-of-data)</p> <p>If data are unavailable, please indicate the exact reasons why data cannot be made available in a suitable public repository or upon request, including any conditions related to ethical approval, consent from study subjects, commercial or legal restrictions, etc.</p> <p>For data that are available under restricted access, the Data Availability statement must specify</p> <ul style="list-style-type: none"> - the reasons for access restrictions - what the restrictions are - how one can get access to the data - who to contact to request access - any restrictions on who the data can be made available to or for which purpose - the expected timeframe for response to access requests - for how long the data will be available once access has been granted. 	<p>Our work satisfied this requirement by providing the behavioral data and processed fMRI data via an OSF link, and it is informed in the Data availability section. No clinical or third party data are included.</p>
<p>Please use the following template to provide all the information stated above:</p> <p>The XX data generated in this study have been deposited in the YY database under accession code ZZ [add hyperlink here]. The XX data are available under restricted access for {insert reason}, access can be obtained by {explain how}. The raw XX data are protected and are not available due to data privacy laws. The processed XX data are available at YY. The XX data generated in this study are provided in the Supplementary Information/Source Data file. The XX data used in this study are available in the YY database under accession code ZZ [Add hyperlink here].</p>	<p>The data availability section follows the guide as below:</p> <p>The human behavioral data used in this study are available in the GitHub repository at https://github.com/brain-machine-intelligence/Rdim-mvpa-model. The processed fMRI data generated in this study (ROI-masked EPI) have been deposited in the OSF database (https://osf.io/2gyue). The data used to create the figures in this paper are provided in the Source Data file. Source data are provided with this paper.</p>
<p>We notice that you have deposited your code in a Github repository, which we fully support. We strongly encourage you in addition to make your code citable by obtaining a DOI for the Github repository in order to provide a permanent reference to the version of the code used in this study and improve reproducibility. This can be done by linking the repository to Zenodo, following the instructions here: https://guides.github.com/activities/citable-code/ Please cite the Github repository in your manuscript text or Code Availability statement and in your reference list: authors, title (this paper), repository name, DOI identifier, year.</p>	<p>We obtained a DOI for the Github repository and cite it in the Code Availability statement and the reference list:</p> <p>The code used for the neural and simulation analyses in this study is available in the GitHub repository (https://github.com/brain-machine-intelligence/Rdim-mvpa-model) and archived at Zenodo (https://doi.org/10.5281/zenodo.17412741).</p>

Alternatively, you can deposit the code in Gigantum or Code Ocean for the same purpose.

Methods

Our guidance:

Your response:

For human studies, indicate the sex and/or gender, number and age of participants in every experiment, and provide a statement on whether informed consent was obtained by participants in the Reporting Summary and Methods. Please also provide information on participant compensation

To comply with SAGER guidelines, the Reporting Summary and Methods should include whether sex and/or gender was considered in the study design and whether sex and/or gender of participants was determined based on self-report or assigned (and methodology used). If no sex or gender analysis was carried out, please clarify why.

Data should be reported disaggregated for sex and gender where this information has been collected and consent has been obtained for reporting and sharing individual-level data; disaggregated numbers for individual experiments must be provided in the source data files whereas overall numbers may be provided in the methods section and Nature Portfolio Reporting Summary.

For more information please see
<https://www.nature.com/articles/s41467-022-30398-1>

Sufficient details of the experiments must be provided in the Methods section such that they could be reproduced without reference to published papers. Use of the term "as described previously" is not encouraged.

The requirements are fulfilled in the Participant section of the manuscript and the Reporting Summary file.

The manuscript updated during the previous round revision to include sufficient methodological details as a stand-alone scientific paper, fully addressing the reviewer's concerns.

End matter

Our guidance:

Your response:

Nature Portfolio defines Competing Interest (CI) as financial and non-financial interests (including but not limited to funding, employment, stocks, shares, patents, personal or professional relationships with individuals or institutions, and unpaid membership advocacy) that could be perceived to directly undermine the objectivity, integrity, and value of a publication, or could be seen as having an influence on the judgments and actions of authors with regard to objective data presentation, analysis, and interpretation.

Please thoroughly review our policy on Competing Interests and include a detailed statement both in your final manuscript file and in our manuscript tracking system. Please ensure the statements are identical in both. Be specific about how each point stated relates to the research and list applicable author initials, and/or patent numbers.

If there are no competing interests, a negative statement must be included.

<https://www.nature.com/nature-research/editorial-policies/competing-interests>

The authors declare no competing interests (statement is identical in the manuscript).

Please confirm that all relevant funding awarded to each author is described in the Acknowledgements section. List each grant number, followed by the initials of the author who received it.

According to all the authors' check, the single grant source was included as below:

[This research was supported by the National Research Foundation of Korea \(NRF\) funded by the Korean government \(MSIT\) \(No. RS-2024-00439903\).](#)

Additional Revisions

Our guidance:

Your response:

<p>For any Supplementary Figures, please check and confirm that:</p> <ul style="list-style-type: none"> * If data is presented as bar charts, individual data points are shown using overlaid dot plots. * The n number (i.e. the sample size used to derive statistics) is provided and defined as a precise value (not a range), using the wording "n=X samples/cells/independent experiments" etc. where applicable. * Any chart axis, error bars, scale bars, molecular weight markers, symbols and colour scales are defined. * Any statistical tests used for data analysis are specified and exact p-values are provided either on the figures themselves, in the legend or in the Source Data file. * Wherever representative data such as blots or micrographs are shown, the legend indicates how many times the experiment was repeated with similar results. * Full uncropped scans of any cropped gel/blot images are provided as an additional Supplementary Figure or in the Source Data file. 	
	All requirements were fulfilled.

Preparing your manuscript files

Our guidance:

Your response:

<p>Unless otherwise stated please limit individual file sizes to approximately 30MB. We strongly encourage the use of repositories for large datasets or source data due to size considerations.</p>	<p>Some panels in Figure 1 contain many data points, making the .ai file >30 MB; alternative acceptable formats are supplied. All other files are <30 MB.</p>
<p>Please supply a brief (maximum 250 characters, including spaces) summary of the main findings of the paper to be used on our website and in our e-alerts. The summary should be written in the third person in language suitable for a broad audience. The summary may be edited by the editors prior to publication. Please provide this summary in your cover letter.</p>	<p>Included in the cover letter: Behavior and fMRI reveal that the human lateral prefrontal cortex factorizes goal and uncertainty into a geometric code, balancing flexible goal pursuit with stable control to prevent erratic behavior.</p>
<p>Please supply the main manuscript file in either Microsoft Word or LaTeX format</p>	<p>The manuscript file is provided in LaTeX format.</p>
<p>Please provide figures as individual vector files with editable text. Acceptable file types for figures are .ai, .eps, .pdf, .ppt or Chem Draw for fully editable vector-based art. For detailed guidance on figure preparation, see https://www.nature.com/documents/aj-artworkguidelines.pdf</p>	<p>The main figures are provided as .ai files with the original .ppt files as well.</p>
<p>Please note that all Supplementary Information must be provided as a single separate PDF file, not within the manuscript file.</p> <p>All Supplementary Information items (e.g. Supplementary Figures, Supplementary Tables, Supplementary Methods, Supplementary Notes, Supplementary Discussion, Supplementary References) must be included in one PDF document. Please refer to our formatting guide when preparing your supplementary information file: https://www.nature.com/documents/ncomms-formatting-instructions.pdf</p>	
<p>All Supplementary Information files (e.g. Supplementary Data, Supplementary Software, etc.) must be cited in the main text.</p> <p>Every Supplementary Figure must be accompanied by a legend of up to 350 words, referring to all panels, and a brief title that summarises the whole figure.</p> <p>Only Supplementary Movie, Audio, Data and Software files should be submitted separately from the Supplementary Information.</p>	<p>All requirements were fulfilled.</p>

<p>The use or adaptation of previously published images is strongly discouraged. If this is unavoidable, please request the necessary rights documentation to re-use such material from the relevant copyright holders and return this to us when you submit your revised manuscript. Please check whether your manuscript or Supplementary Information contain third-party images, such as figures from the literature, stock photos, clip art or commercial satellite and map data.</p> <p>If any elements of your submitted work have been created with BioRender you will need to ensure you have obtained a publication license from BioRender, adhering to the user requirements as outlined within the license. The reference for BioRender created graphics should be present in the accompanying legend of the display material it is present in.</p> <p>A copy of the publication license should be uploaded to our system as a related manuscript file upon resubmission.</p> <p>For more information please see the BioRender knowledge article here: https://help.biorender.com/hc/en-gb/articles/21283116932765-CC-BY-publishing-and-reader-permissions</p> <p>For more information on what constitutes ownership by a third party, please contact our Editorial Assistant at naturecommunications@nature.com</p> <p>Please check in particular:</p> <p>Please note that suspected third party content is present in figures 1a,b,e; 3a and Supplementary Figure 2; Supplementary Figure 6a; Supplementary Figure 9a</p>	<p>We carefully reviewed all figures for potential third-party content. Panels 1a,b,e; S6a; S9a are newly redrawn from scratch and are not adapted from any published images. Figure 3a and Supplementary Figure 2 are our original work.</p>
--	---

Forms to complete

Our guidance:

Reporting Summary

Please revise the Reporting Summary according to the requests below. After making the requested changes, please be sure to include the final version of your Reporting Summary in your submission as a supplementary information file. Please note that this form is a dynamic 'smart pdf' and must therefore be downloaded and completed in Adobe Reader, instead of opening it in a web browser.

Please update your current checklist or download from:

<https://www.nature.com/documents/nr-reporting-summary.pdf>

Your response:

The Reporting Summary has been updated per instructions.

Reporting Summary

Our guidance:	Your response:
Software	
Please ensure all the data collection/data analysis software/tools/algorithms/packages used in the study are clearly mentioned in the manuscript and are also listed in the reporting summary (with version numbers).	We clearly mentioned all the required details.
Data Availability	
Please provide a complete data availability statement in the manuscript and in the reporting summary.	The complete statement was provided under the Data section.
Field Specific Reporting	
Life Sciences Study Design	
Please describe the general measures taken to verify the reproducibility of the experimental findings. If all attempts at replication were successful, confirm this OR if there are any	We verified reproducibility by openly sharing all underlying data and analysis code and providing a consolidated Source Data file for every figure. All key results were

findings that were not replicated or cannot be reproduced, note this and describe why.	reproducible with the shared pipeline; we are not aware of any findings that failed to replicate.
--	---

You will need to upload:

Completed Third Party Rights Table (if relevant)	Not applicable.
A point-by-point response to the reviewers' comments	All issues have been addressed. Please refer to the Author Checklist under "Remaining Reviewer Comments" for details. (The same contents in the 'Rebuttal letter.pdf')
A completed copy of this checklist	Author_checklist.docx
The main manuscript file in either Microsoft Word or LaTeX format	A LaTeX format manuscript file (main.tex)
Separate Figure files	Figure1 (.ai, .pptx, .svg), Figure2 (.ai, .pptx), Figure3 (.ai, .pptx), Figure4 (.ai, .pptx),
Separate Source Data files	Source Data.xlsx
Inventory of Supporting Information	si.pdf

All issues have been addressed. Please refer to the Author Checklist under "Remaining Reviewer Comments" for details.

ARTICLE IN PRESS

¹ Factorized embedding of goal and uncertainty in the
² lateral prefrontal cortex guides stably flexible learning

³ Yoondo Sung¹, Mattia Rigotti², and Sang Wan Lee^{1, 3, 4, 5, 6, *}

⁴ ¹Department of Bio and Brain Engineering, Korea Advanced Institute of Science and Technology
⁵ (KAIST), Daejeon, Republic of Korea.

⁶ ²IBM Research, Zurich, Switzerland.

⁷ ³Center for Neuroscience-inspired AI, Korea Advanced Institute of Science and Technology
⁸ (KAIST), Daejeon, Republic of Korea.

⁹ ⁴Department of Brain and Cognitive Sciences, Korea Advanced Institute of Science and
¹⁰ Technology (KAIST), Daejeon, Republic of Korea.

¹¹ ⁵Kim Jaechul Graduate School of AI, Korea Advanced Institute of Science and Technology
¹² (KAIST), Daejeon, Republic of Korea.

¹³ ⁶Graduate School of Data Science, Korea Advanced Institute of Science and Technology (KAIST),
¹⁴ Daejeon, Republic of Korea.

¹⁵ *Correspondence: sangwan@kaist.ac.kr (S.W.L.)

16 Abstract

17 A major challenge for adaptive agents is achieving behavioral flexibility without compromis-
18 ing stability—particularly in goal-directed learning within uncertain environments. Agents
19 must adjust as goals shift while maintaining resilience against noisy signals, necessitating
20 the delicate tradeoff: balancing flexibility for goal pursuit with stability for preventing er-
21 ratic behavior. To investigate how the brain navigates this dilemma, we combined model
22 simulations with behavioral and fMRI data collected during a goal-directed learning task
23 under varying levels of uncertainty. Our simulations revealed that model-free learning strug-
24 gles with the flexibility-stability trade-off, whereas model-based learning allows for flexible
25 goal pursuit with varying degrees of stability. Interestingly, human participants displayed
26 both stable and flexible goal-directed behavior. The fMRI data uncovered the underlying
27 mechanism: goals and uncertainty are represented as factorized embeddings in the lateral
28 prefrontal and orbitofrontal cortex. Notably, the neural separability of goals and their re-
29 silience to uncertainty in these regions correlated with participants' behavioral flexibility and
30 stability.

31 Introduction

32 Biological agents must rapidly adapt learned behavior, the ability known as behavioral flex-
33 ility. As a core feature of adaptive behavior, goal-directed learning requires representing
34 action–outcome contingencies and updating them to attain a goal. Under uncertainty, how-
35 ever, it becomes challenging to distinguish goal-relevant feedback from environmental noise.
36 Being overly sensitive to noisy outcomes can undermine stability, whereas failing to revise
37 outdated representations hinders adaptation to changing contexts. To illustrate the stabil-
38 ility–flexibility dilemma, consider a rescue operation under highly uncertain conditions. A
39 rescue team must remain flexible enough to adjust if new information pinpoints the sur-
40 vivors' location or the building's accessibility. However, overreacting to every rumor wastes
41 resources on false leads. In such a scenario, filtering out noisy events (stability) while re-
42 sponding appropriately to significant signals (flexibility) is crucial for success under shifting
43 circumstances (stably flexible decision-making).

44 The prefrontal cortex (PFC) plays a vital role in cognitive control, which is necessary for
45 behavioral flexibility^{1,2}. It specializes in representing and integrating various variables such
46 as context, decisions, and sensory information^{3,4,5,6}. The PFC guides context-dependent
47 decision-making by flexibly altering its information representation based on the current
48 goals^{3,7,8,9,10}. Neural representations in the PFC reflect task demands, facilitating successful
49 goal achievement and directly impacting performance^{11,6}.

50 The PFC also encodes environmental uncertainty. The lateral prefrontal cortex (LPFC)
51 and orbitofrontal cortex (OFC) are particularly important in representing and tracking en-
52 vironmental uncertainty to adjust behavior^{12,13,14,15}. The LPFC is involved in arbitrating
53 between different learning strategies according to uncertainty^{16,17,18,19}. It is also known to
54 perform an important role in model-based control necessary for goal pursuit²⁰. The OFC
55 is involved in making confidence judgments under uncertainty^{21,22}, representing state-space

56 according to task demands^{23,24,25}, and calculating values^{26,27,28,29}.

57 While we must flexibly adjust our behavior to pursue goals, we also need to respond
58 stably to changes in an uncertain environment. However, cognitive stability and flexibility
59 are generally considered to be in a trade-off relationship, making it challenging to achieve
60 both^{30,31,32,33,34,35}. Similarly, neural separability for flexibility and neural robustness for
61 stability appear to be conflicting properties. In situations where multiple pieces of infor-
62 mation are mixed, high-dimensional neural representations that easily distinguish stimuli
63 or contexts are advantageous for flexible responses^{11,36,37,38,39}. Conversely, stable task per-
64 formance requires neural representations that abstract information into a low-dimensional
65 manifold by removing task-irrelevant noise or distractors^{40,41,42,7}. This creates a known
66 trade-off between neural separability and robustness, depending on the representational di-
67 mensionality^{37,43,44,42}. While these studies focused on perceptual decision-making to explore
68 neural representations for cognitive flexibility or stability, it remains unclear how the brain
69 achieves stable, flexible goal pursuit during sequential decision-making with various contex-
70 tual changes.

71 Our study employs model simulations, behavioral data analysis, and model-based fMRI
72 analysis, to examine whether and how the brain resolves the behavioral flexibility-stability
73 trade-off during goal-directed learning. We compare humans' behavioral flexibility and sta-
74 bility with conventional value learning theory, including model-free and model-based rein-
75 forcement learning. We then analyze how the neural representations of goals and uncertainty
76 in related brain areas, including the dlPFC, vLPFC, and OFC, facilitate stable and flexible
77 goal pursuit.

78 **Results**

79 **Behavioral stability-flexibility dilemma during goal-directed learning**

80 To explore the neural representations involved in pursuing goals within uncertainty-changing
81 environments, we used behavior and fMRI data published in a previous study¹⁷. Twenty par-
82 ticipants performed a context-dependent two-stage Markov decision task (Fig. 1a-b; Task).
83 Each state in the task was represented by a fractal image, and participants made left or
84 right choices at each state across two stages. After the second choice, an outcome state
85 and a corresponding coin were displayed, with the task's objective being to maximize the
86 accumulated coin score.

87 The experiment featured four types of block conditions, corresponding to each combina-
88 tion of two binary context conditions: *goal condition* and *uncertainty condition* (Fig. 1b).
89 By systematically varying goal specificity and state transition uncertainty, respectively, this
90 design captures two important dimensions of goal pursuit: the clarity of the goal to pursue
91 and the predictability of the environment. First, the *goal condition* includes *specific goal* and
92 *non-specific goal* conditions. In the *specific goal condition*, only coins matching the color (red,
93 blue, or yellow) of a box presented at the start of the trial were considered valid rewards,
94 and coins of other colors scored zero. For example, when a red box appears, participants
95 understand that the task for that trial is to reach the state yielding a red coin; the goal is
96 therefore the red-coin state. This setup encourages participants to align their actions with
97 a clear goal, thereby promoting goal-directed behavior. Conceptually, it resembles having a
98 specific goal (e.g. rescuing survivors in a single building). In contrast, the *non-specific goal*
99 *condition* served as a control, where a white box was displayed, and coins of any color were
100 recognized as valid rewards. While this condition still involves reward maximization, the
101 absence of a clearly defined goal can lead to more habitual behaviors that are less sensitive
102 to environmental changes. It is comparable to having a broad, general objective (e.g. “help

¹⁰³ anyone, anywhere”).

¹⁰⁴ The *uncertainty condition* modulated the distribution of state-action-state transition
¹⁰⁵ probabilities. Under the *low uncertainty condition*, the probabilities of the two possible
¹⁰⁶ subsequent states following a choice were assigned as (0.9, 0.1), whereas in the *high uncer-*
¹⁰⁷ *tainty condition*, these probabilities were (0.5, 0.5). This difference captures how predictable
¹⁰⁸ (or unpredictable) an environment is. In the low uncertainty condition, state transitions
¹⁰⁹ follow relatively predictable probabilities, whereas in the high uncertainty condition, transi-
¹¹⁰ tions are dominated by random variability. Such uncertainty changes influence the efficacy of
¹¹¹ goal-directed learning, since increased uncertainty reduces the prediction accuracy of state-
¹¹² action-state transitions, a key factor in model-based learning.

¹¹³ We investigated the impact of goal and uncertainty contexts on state-action values and
¹¹⁴ behavior. Using an oracle agent, we calculated the true state-action values for each context
¹¹⁵ using the Bellman equation⁴⁵. The results showed that in the *low uncertainty condition*,
¹¹⁶ where the state transition probabilities were biased, the value difference between left and
¹¹⁷ right choices was larger than in the high uncertainty condition (Fig. 1c). From a value-
¹¹⁸ based decision-making perspective, in the *high uncertainty condition*, participants are less
¹¹⁹ likely to make optimal choices due to smaller value difference^{46,47,48}. However, by measuring
¹²⁰ *choice optimality* and *choice consistency* for each trial (Behavioral measures), we found
¹²¹ uncertainty did not affect task performance in the *specific-goal condition* (Fig. 1d). In the
¹²² *non-specific goal condition*, high uncertainty led to fewer optimal choices (Fig. 1d, average
¹²³ choice optimality; low unc.: 0.850(mean) \pm 0.035(s.e.m.); high unc.: 0.619 \pm 0.021; low unc.
¹²⁴ vs high unc.: $t(19) = 5.975, p < 0.0001$) and more frequent choice changes (Fig. 1d, average
¹²⁵ choice consistency; low unc.: 0.897 ± 0.023 ; high unc.: 0.847 ± 0.021 ; low unc. vs high unc.:
¹²⁶ $t(19) = 5.294, p < 0.0001$). Surprisingly, in the *specific goal condition*, choice optimality and
¹²⁷ consistency were not significantly affected by high uncertainty (Fig. 1d; low unc. vs high unc.
¹²⁸ (optimality): $t(19) = 0.887, p = 0.384$; low unc. vs high unc. (consistency): $t(19) = 0.891$,

¹²⁹ $p = 0.386$), suggesting that humans achieve robust performance when pursuing a specific
¹³⁰ goal in uncertain conditions. Relatedly, participants' response times were longer in the
¹³¹ specific-goal than in the non-specific-goal condition, and within the specific-goal condition
¹³² they increased further under high uncertainty (Supplementary Fig. 1). These prolonged
¹³³ latencies likely reflect the additional deliberation required both to reach the specified goal
¹³⁴ and to handle unexpected state transitions in an uncertain environment.

¹³⁵ This finding led us to focus on the *specific goal condition* to understand humans' stably
¹³⁶ flexible goal pursuit. Optimal performance requires flexible behavior aligned with the goal,
¹³⁷ unhindered by noisy state transitions. We quantified these attributes using *choice versatil-*
¹³⁸ *ity* and *choice consistency* (Fig. 1e, Behavioral measures). Each measure, calculated on a
¹³⁹ trial-by-trial basis, was designed to quantify choice sensitivity to goal changes and choice
¹⁴⁰ consistency despite environmental noise. The session-wide average of these trial-level values
¹⁴¹ provided an individual measure of each participant's behavioral flexibility or stability. In the
¹⁴² *specific goal condition*, participants must select different trajectories to reach the outcome
¹⁴³ state matching the given goal (red, blue, or yellow). Hence, when the goal of the current
¹⁴⁴ trial is different from the previous trial, a different choice at the same state indicates choice
¹⁴⁵ versatility = 1, whereas the same choice indicates choice versatility = 0. Conversely, given
¹⁴⁶ a goal, participants should consistently make the optimal choice despite noisy state transi-
¹⁴⁷ tions. Therefore, if the goal of the current trial matches the goal of the previous trial, making
¹⁴⁸ the same choice at the same state indicates choice consistency = 1, while a different choice
¹⁴⁹ indicates choice consistency = 0.

¹⁵⁰ To understand how human participants maintain the behavioral flexibility-stability bal-
¹⁵¹ ance, we compared the relationships between flexibility, stability, and performance of human
¹⁵² participants with those of standard value learning models, including model-based (MB)
¹⁵³ and model-free (MF) reinforcement learning algorithms (Fig. 1f, Simulation). The over-
¹⁵⁴ all performance of each participant was calculated by averaging trial-wise *choice optimal-*

ity (Behavioral measures), which evaluated the degree to which their actions aligned with those of an oracle agent possessing complete knowledge of reward contingencies and state-transition probabilities (1 = optimal, 0 = non-optimal). Because the MF algorithm learns a reward function based on consistent sampling of action-outcome events, whereas the MB algorithm does so by learning the dynamics of the action-state-outcome, each model's behavior is known to be stable and flexible, respectively. We found effects of flexibility and stability on performance in MB agents and humans, but not in MF agents (the left and the middle plot of Fig. 1f). For the flexibility–performance comparison, humans showed a strong positive correlation ($r = 0.827$, $t(18) = 6.244$, $p < 0.0001$), MB agents showed a comparable effect ($r = 0.799$; 95% subsampling CI [0.649, 0.916], 10000 resamples with subsample size $m = 20$), whereas MF agents showed no reliable relationship ($r = -0.035$; CI [-0.475, 0.411]). For the stability–performance comparison, humans again displayed a robust association ($r = 0.779$, $t(18) = 5.269$, $p < 0.0001$), MB agents exhibited an even stronger link ($r = 0.963$; CI [0.933, 0.986]), while MF agents remained near zero ($r = 0.042$; CI [-0.399, 0.485]). Furthermore, MF agents suffer from a flexibility-stability trade-off (as indicated by a strong negative correlation in the right plot of Fig. 1f), whereas humans and the MB agents did not. Humans showed a moderate positive flexibility–stability relation ($r = 0.541$, $t(18) = 2.726$, $p = 0.014$), MB agents showed a strong positive coupling ($r = 0.829$; CI [0.681, 0.930]), and MF agents showed a pronounced negative coupling ($r = -0.964$; CI [-0.986, -0.929]). Notably, humans exhibited the most flexible and stable behavior.

¹⁷⁶ Evidence of goal and uncertainty representation in PFC during goal-
¹⁷⁷ directed learning

¹⁷⁸ As the first step to investigate neural underpinnings of flexible and stable goal pursuit,
¹⁷⁹ we investigated whether the brain encodes goals and uncertainty information. We used
¹⁸⁰ the ROI-based multivoxel pattern analysis, focusing on eight brain regions known to be
¹⁸¹ engaged in context-dependent behavior: ventrolateral prefrontal cortex (vlPFC), dorsolateral
¹⁸² prefrontal cortex (dlPFC), orbitofrontal cortex (OFC), anterior cingulate cortex (ACC), pre-
¹⁸³ supplementary motor area (preSMA), primary visual cortex (V1), hippocampus (HPC), and
¹⁸⁴ ventral striatum (vStr). The vlPFC is known to be involved in the arbitration of MB/MF
¹⁸⁵ reinforcement learning^{17,49,50,18}, whereas the dlPFC is known to represent context information
¹⁸⁶ and guide task-switching^{51,52,15}. The OFC is implicated in state-space representation^{23,24,25}.
¹⁸⁷ The ACC plays a crucial role in monitoring conflicts in guiding adaptive adjustments in
¹⁸⁸ cognitive control^{53,54}. The preSMA is involved in the flexible control of voluntary actions⁵⁵.
¹⁸⁹ The HPC is involved in the formation of cognitive maps^{56,57,58}. Lastly, the ventral striatum
¹⁹⁰ plays a key role in computing reward prediction errors^{59,60,61}. The AAL3 atlas⁶² was used
¹⁹¹ to define each ROI (Supplementary Fig. 2).

¹⁹² For flexible goal pursuit, the brain must effectively encode information about goals and
¹⁹³ associated states. We focused on the specific goal condition and conducted a decoding
¹⁹⁴ analysis on multivoxel patterns to quantify goal-related information. This was assessed by
¹⁹⁵ the classification accuracy of three distinct goals (red, blue, and yellow) using linear support
¹⁹⁶ vector machines (SVM) (fMRI decoding analyses). All classifications on the fMRI data were
¹⁹⁷ performed with leave-one-run-out cross-validation. To prevent potentially biased predictions
¹⁹⁸ from class imbalance, we performed 100 rounds of undersampling to balance the classes,
¹⁹⁹ using the average accuracy from these iterations as the final result. After this calibration,
²⁰⁰ we confirmed that the decoding accuracy did not significantly exceed the chance level when

201 the class label is shuffled (Supplementary Fig. 3). We trained and tested separate SVM
 202 classifiers for each trial event (fix: fixation, S1: state 1, A1: action 1, S2: state 2, A2:
 203 action 2, S3: state 3; Fig. 1a) to obtain decoding accuracy specific to that event. The
 204 decoding accuracy was calculated as a single value per participant based on multiple trials.
 205 Group-level statistical tests were performed, treating each participant as a random sample.

206 The vlPFC, dlPFC, OFC, ACC, preSMA and V1 demonstrated the highest goal decoding
 207 accuracy during the second stage (Fig. 2a; event-specific decoding accuracies; see Supple-
 208 mentary Table 2 for the statistical details). The average decoding accuracy across all events
 209 corroborates that these six regions significantly represented goal information (Fig. 2b; t-test
 210 of event-averaged decoding accuracy against the chance level (33.3%); vlPFC: $37.2\% \pm 0.7\%$,
 211 $t(19) = 5.712$, $p < 0.0001$; dlPFC: $38.8\% \pm 0.9\%$, $t(19) = 5.999$, $p < 0.0001$; OFC:
 212 $36.6\% \pm 0.7\%$, $t(19) = 4.443$, $p = 0.0003$; ACC: $35.4\% \pm 0.6\%$, $t(19) = 3.459$, $p = 0.0026$;
 213 preSMA: $35.4\% \pm 0.6\%$, $t(19) = 3.701$, $p = 0.0015$; V1: $36.9\% \pm 0.5\%$, $t(19) = 7.401$, $p <$
 214 0.0001 ; uncorrected per predefined ROI; Supplementary Table 2). Furthermore, regarding
 215 state information, the vlPFC, dlPFC, OFC, ACC, preSMA, V1, and HPC demonstrated
 216 significant decodability of both intermediate and outcome states (Supplementary Fig. 4a-b).
 217 Notably, decoding accuracy for the intermediate state peaked during stage 2, while outcome
 218 state decoding peaked during stage 3.

219 For stable and flexible goal pursuit, the neural representation of goals must be separable
 220 from uncertainty. This demands uncertainty encoding during goal-directed learning. When
 221 measured as the test accuracy of classifying uncertainty conditions using a linear SVM on
 222 multivoxel patterns in the specific goal condition and the non-specific goal conditions, the
 223 vlPFC, dlPFC, and OFC significantly represented uncertainty information exclusively within
 224 the specific goal condition (Fig. 2d; t-test of event-averaged decoding accuracy against the
 225 chance level (50%); vlPFC: $51.8\% \pm 0.7\%$, $t(19) = 2.549$, $p = 0.020$; dlPFC: $52.1\% \pm 0.9\%$,
 226 $t(19) = 2.412$, $p = 0.026$; OFC: $53.4\% \pm 0.9\%$, $t(19) = 3.749$, $p = 0.0014$; uncorrected per

²²⁷ predefined ROI; Supplementary Table 2). Additionally, our principal component analysis
²²⁸ (PCA) showed that the neural dimensionality in vlPFC, dlPFC, OFC, and ACC is higher
²²⁹ in the specific goal condition compared to the non-specific goal condition (Supplementary
²³⁰ Fig. 5). These results imply that the LPFC and OFC encode uncertainty while engaging in
²³¹ complex neural computations to guide goal-directed behavior.

²³² Factorized embedding of goal and uncertainty in the LPFC

²³³ Building on our findings that goal and uncertainty are represented in the LPFC and OFC,
²³⁴ we sought to investigate how these two variables are represented in a single neural space
²³⁵ to facilitate flexible yet stable goal pursuit. Following previous studies on representational
²³⁶ geometry^{11,41,43,63}, we evaluated three hypotheses on mixed representations of goals and un-
²³⁷ certainty (Fig. 3a). The types of possible linear separations vary depending on the complexity
²³⁸ of the neural embedding structure^{64,65}.

²³⁹ If only one of the two variables is represented, one can linearly separate the neural
²⁴⁰ representations of the classes of that variable, but not of the other variable (compression
²⁴¹ hypothesis; the first column of Fig. 3a). The single represented variable remains stable and
²⁴² invariant across the changes in the other variable. However, because the represented variable
²⁴³ has no information about the other variable, it is impossible to distinguish situations in which
²⁴⁴ that other variable changes. For instance, if the goal is represented without representing
²⁴⁵ uncertainty, the distinct goal information remains accessible but cannot detect or adapt to
²⁴⁶ fluctuations in uncertainty.

²⁴⁷ Conversely, when both variables are independently represented along their respective cod-
²⁴⁸ ing axes, binary classifications involving the two variables can be linearly separated (factor-
²⁴⁹ ized mixing hypothesis; the second column of Fig. 3a). In this case, both variables maintain
²⁵⁰ distinct representations, ensuring that changes in one do not alter the embedding structure of
²⁵¹ the other. As a result, a downstream neural readout that decodes one variable can generalize

252 across variations in the other variable. In such a factorized embedding structure, it is pos-
 253 sible to detect uncertainty changes, while goal information remains consistently represented
 254 regardless of uncertainty levels.

255 Lastly, if there exists an interaction between the two variables, where the coding axis of
 256 one variable changes contingent on the other, dichotomies involving the nonlinear interaction
 257 can also be linearly separated (nonlinear mixing hypothesis; the third column of Fig. 3a).
 258 This high-dimensional neural embedding structure allows distinction across a wide range of
 259 situations arising from variable combinations. However, such dependency among variables
 260 reduces generalizability. As illustrated in the figure, uncertainty changes lead to changes in
 261 the goal embedding structure, making the representation of goal information highly sensitive
 262 and vulnerable to uncertainty shifts.

263 To understand the representational geometry, we performed a *shattering analysis* (Shat-
 264 tering analysis) to identify the types of possible linear separations among all dichotomies. We
 265 categorized all dichotomies into four types: goal, uncertainty, linear, and nonlinear (Supple-
 266 mentary Fig. 6). The average test accuracy of dichotomies within each category was defined
 267 as the shattering dimensionality (SD) for that category. Figure 3a illustrates plausible neural
 268 embeddings under different combinations of these four types of separability. If the neural
 269 embedding follows the compression hypothesis, one of the SDs for goal or uncertainty will
 270 be distinctly high. Under the factorized mixing hypothesis, the SDs for goal, uncertainty,
 271 and linear will be notably high, whereas the SD for nonlinear will be significantly lower.
 272 Conversely, if goal and uncertainty form a nonlinearly mixed embedding, all four types of
 273 SD will be substantially high.

274 To evaluate which dichotomies—averaged into four SD categories—were linearly sepa-
 275 rable within each ROI, we applied multivoxel-pattern linear decoding. Aside from incor-
 276 porating multiple binary label sets, the procedure followed our standard decoding pipeline
 277 (fMRI decoding analyses). For each ROI, we trained separate linear SVM classifiers at every

278 task-informative event (S1, A1, S2, A2, S3, fix') and averaged their decoding accuracies to
 279 obtain the SD score. The initial fixation period (fix) was excluded because it lacks task-
 280 relevant information for the current trial. However, as the fixation period preceding a new
 281 trial retains residual information from the previous trial, neural activity for the current trial
 282 was decoded using the fixation epoch of the subsequent trial (fix').

283 The vlPFC, dlPFC, and OFC showed a separability profile corresponding to factorized
 284 mixing (Fig. 3b). These brain regions showed significant SD for goal, uncertainty, linear,
 285 and nonlinear types (Supplementary Fig. 7; vlPFC goal: 0.530 ± 0.005 ($t(19) = 5.634$, $p <$
 286 0.0001); vlPFC uncertainty: 0.518 ± 0.007 ($t(19) = 2.582$, $p = 0.0183$); vlPFC linear:
 287 0.519 ± 0.003 ($t(19) = 5.527$, $p < 0.0001$); vlPFC nonlinear: 0.508 ± 0.002 ($t(19) = 4.018$, $p =$
 288 0.000735). dlPFC goal: 0.541 ± 0.007 ($t(19) = 6.259$, $p < 0.0001$); dlPFC uncertainty:
 289 0.521 ± 0.009 ($t(19) = 2.397$, $p = 0.0270$); dlPFC linear: 0.526 ± 0.004 ($t(19) = 6.158$, $p <$
 290 0.0001); dlPFC nonlinear: 0.511 ± 0.002 ($t(19) = 4.432$, $p = 0.000286$). OFC goal: $0.527 \pm$
 291 0.006 ($t(19) = 4.663$, $p = 0.00017$); OFC uncertainty: 0.534 ± 0.009 ($t(19) = 3.724$, $p =$
 292 0.00144); OFC linear: 0.521 ± 0.004 ($t(19) = 4.807$, $p = 0.000122$); OFC nonlinear: $0.507 \pm$
 293 0.003 ($t(19) = 2.904$, $p = 0.00909$); paired t -tests against the chance level (0.5) within
 294 each predefined ROI, uncorrected). Notably, the vlPFC and dlPFC showed significantly
 295 lower nonlinear SD compared to goal and linear SDs (pairwise comparison of event-averaged
 296 SDs between four classification types by paired t -test; see Supplementary Table 3 for the
 297 statistical details). Similarly, the OFC showed significantly lower nonlinear SD than goal,
 298 uncertainty, and linear SDs. According to our three hypotheses (Fig. 3a), these results
 299 demonstrated a factorized embedding of goal and uncertainty. It also suggests a neural
 300 mechanism that maintains stable goal representations across varying levels of uncertainty.

301 The next step was to clarify how factorized neural embeddings contribute to maintaining
 302 stable and flexible goal-pursuit behavior. To this end, we conducted correlation analyses
 303 quantifying the relationship between different patterns of neural separability and task be-

³⁰⁴ havior. The neural metric was restricted to signals that directly inform the choice behavior.
³⁰⁵ Accordingly, outcome-related events occurring after the choice (S3 and fix') were excluded,
³⁰⁶ and each SD was recalculated as the mean accuracy across the remaining pre-outcome events
³⁰⁷ (S1–A2). Multiple comparison corrections were applied with a false-discovery rate (Ben-
³⁰⁸ jamini–Hochberg procedure) for the number of ROIs with $q = 0.05$. Adjusted p-values were
³⁰⁹ reported for all exploratory correlation analyses (Statistical analysis).

³¹⁰ The significant correlation between each type of neural separability and behavioral mea-
³¹¹ sures that we observed suggests that goal separability in the vLPFC, dlPFC, and OFC is
³¹² associated with goal pursuit performance (Fig. 3c; FDR-corrected, $q < 0.05$; Statistical
³¹³ analysis). Notably, higher goal separability in the LPFC correlates with greater behavioral
³¹⁴ flexibility, stability, and performance. These results carry two critical implications. First,
³¹⁵ a neural embedding in the LPFC, distinguishing goals from uncertainty, is associated with
³¹⁶ goal-dependent behavioral adaptation. Second, a clear goal representation is essential for
³¹⁷ consistently pursuing desired outcomes across multiple stages, even in noisy environments.
³¹⁸ On the other hand, uncertainty separability was not significantly related to goal-directed
³¹⁹ behavior, while linear and nonlinear interaction separabilities showed relatively weaker but
³²⁰ generally consistent results with goal separability (Fig. 3d). To summarize, a neural embed-
³²¹ ding capable of clearly distinguishing goals is important for effective goal pursuit.

³²² **Neurally stable goal embedding in LPFC guides stably flexible learn-** ³²³ **ing**

³²⁴ We further hypothesized that goal embeddings remain unaffected by uncertainty to maintain
³²⁵ stable behavior. To test this hypothesis, we employed the neural metric called cross-condition
³²⁶ generalization performance (CCGP; Cross-condition generalization performance)⁴¹, defined
³²⁷ as the generalized accuracy of a linear decoder across different conditions. Specifically, we

328 trained a linear classifier (SVM) to decode a target variable (i.e., the goal) in one context
 329 condition (e.g., low uncertainty) and then tested it in a different condition (e.g., high uncer-
 330 tainty). The CCGP score is calculated as the average test accuracy across these conditions.
 331 If the target embedding is influenced by context, the decision boundary for decoding would
 332 shift, leading to a lower CCGP. Thus, CCGP reflects how stably the neural representation
 333 can be decoded by downstream neural readouts across varying context conditions. To assess
 334 the robustness of goal representations, we measured the goal CCGP across uncertainty lev-
 335 els. Consistent with the preceding decoding analysis, we trained linear SVM classifiers for
 336 each ROI at every task-informative event (S1–fix') and averaged their decoding accuracies
 337 to obtain the final CCGP score.

338 Supporting our hypothesis, the CCGP and SD values in each uncertainty condition were
 339 comparable (Fig. 4a; pairwise comparison of event-averaged neural measures by paired t-
 340 tests within each predefined ROI, uncorrected; see Supplementary Table 4 for the statistical
 341 details). Although the vLPFC, dlPFC, and OFC significantly represented uncertainty in-
 342 formation (Fig. 2b), goal embeddings in these regions remained robust and unaffected by
 343 uncertainty. Similarly, the ACC, preSMA, and V1, which did not significantly represent
 344 uncertainty, maintained uncertainty-robust goal embeddings. The HPC and vStr, the re-
 345 gions that did not significantly represent goals or uncertainty, showed CCGP values near the
 346 chance level.

347 Furthermore, we found that neural robustness in the vLPFC and dlPFC is crucial for goal
 348 pursuit behavior (Fig. 4b; FDR-corrected, $q < 0.05$). The goal CCGPs across uncertainty in
 349 the vLPFC and dlPFC were significantly correlated with behavioral flexibility, stability, and
 350 performance. Additionally, higher CCGP in the OFC was associated with greater behavioral
 351 stability. As in the SD–behavior correlation analysis, we used the mean CCGP averaged
 352 across the pre-outcome events (S1–A2) for this correlation.

353 Additionally, we quantified whether the neural goal representation remained stable across

354 uncertainty levels by computing the parallelism score⁴¹. For each ROI, we averaged the mul-
355 tidimensional BOLD pattern within each class and derived vectors representing the direction
356 from one specific goal representation to another under each uncertainty condition (Supple-
357 mentary Fig. 9a). The cosine similarity between the goal-encoding vectors obtained under
358 the two uncertainty conditions was then measured. Parallelism scores were averaged across
359 five fMRI runs. As with the other classification-based analyses, class balancing was per-
360 formed through undersampling.

361 Consequently, we found that the vlPFC, dlPFC, OFC, and preSMA exhibited signifi-
362 cantly positive parallelism scores (Supplementary Fig. 9b), indicating minimal reorienta-
363 tion of neural goal embeddings in response to uncertainty. This suggests that the underly-
364 ing representational geometry remained largely parallel. Combined with the CCGP results
365 demonstrating robust trial-by-trial decodability based on decision boundaries, these findings
366 confirmed that goal embeddings remain relatively stable under changing uncertainty.

367 Discussion

368 To investigate how humans achieve flexible and stable goal pursuit, we studied human goal-
 369 directed learning in uncertainty-changing environments using a two-stage Markov decision
 370 task. Our findings revealed that humans exhibit more robust behavior during goal-directed
 371 learning. Furthermore, higher behavioral flexibility in response to goal changes correlates
 372 with more stable behavior under uncertainty. We measured the neural representational sepa-
 373 rability and robustness of brain regions, which indicated that goal and uncertainty represen-
 374 tations form factorized embeddings in the vIPFC, dlPFC, and OFC. Neural goal separability
 375 and robustness in these regions are associated with stably flexible goal-directed behavior in
 376 humans.

377 Our study presents a significant conceptual advance by offering a representational solution
 378 to the stability-flexibility dilemma in goal-directed learning. While prior works have primar-
 379 ily focused on identifying ‘what’ variables are encoded in the PFC (e.g., value, uncertainty),
 380 our work elucidates ‘how’ the brain navigates this trade-off through the specific geometrical
 381 structure of its neural representations. Earlier studies in decision-making have attempted to
 382 probe behavioral flexibility in response to changes in goals or context^{66,60,29,67,68,69}, while plac-
 383 ing less emphasis on the concurrent challenge of maintaining stability against environmental
 384 noise. Additionally, studies on uncertainty representation have often centered on perceptual
 385 judgments^{70,71,72,73,74,75} rather than on the complex dynamics of sequential action toward
 386 delayed goals. Our findings bridge this gap by showing, through the lens of representational
 387 geometry, that the PFC employs a joint, factorized coding scheme for goals and uncertainty.
 388 This neural architecture offers a mechanistic account of how the brain reconciles competing
 389 cognitive demands, enabling robust and generalizable goal pursuit while remaining attuned
 390 to environmental statistics, thus supporting both adaptive behavior and resilience to noise.

391 Across the eight ROIs, our MVPA results reveal complementary contributions to stable

yet flexible goal pursuit. V1 encodes task-relevant sensory features, whereas vlPFC, dlPFC, and OFC jointly represent factorized goal and uncertainty information; the strength of these codes correlates with individual differences in behavioral flexibility under uncertainty. By contrast, ACC and preSMA mainly encode goal information, and the magnitude of these signals tracks behavioral performance, suggesting that medial frontal areas relay resolved goals to downstream control systems once uncertainty has been represented in the lateral PFC and OFC circuits^{76,55}. Uncertainty was not decodable in ACC, indicating that our noisy state transitions did not evoke the internal value conflict typically associated with this region^{53,77}, which is consistent with our findings that factorized goal and uncertainty codes help resolve the stability–flexibility dilemma. Hippocampus and ventral striatum showed no reliable pattern-level goal or uncertainty coding within the sensitivity limits of our analysis. Together, these findings support a model in which lateral PFC and OFC furnish a flexible state representation, while medial PFC contributes goal-driven control, collectively balancing behavioral flexibility and stability.

In our study, we confirmed that humans can maintain successful goal pursuit behavior even in uncertain environments. Previous studies have reported that an inaccurate prediction by the model-based learning system¹⁶ or low outcome controllability^{78,79} can reduce goal-directed behavior. High state-transition uncertainty, for example, increases state prediction error, making it more challenging to achieve desirable outcomes. However, our two-stage Markov decision task, which features multiple possible trajectories leading to the goal state across 16 different branches, allows a greater flexibility in reaching the desired outcome. This design contrasts with simpler decision-making tasks, allowing the observation of robust goal-directed behavior under uncertainty. Our findings suggest that this robustness is supported by the factorized representation of goal and uncertainty in the prefrontal cortex, particularly in the vlPFC, dlPFC, and OFC.

The finding in our study that goal representation remains robust to uncertainty is con-

418 sistent with previous findings suggesting that having a factorized representation facilitates
 419 generalization across various contexts^{37,41,80,81,63,82}. Moreover, the perspective that a simi-
 420 lar orthogonal representation structure enables humans to avoid catastrophic forgetting and
 421 learn various tasks^{83,7,84} aligns with the results of our study, indicating that goal embedding
 422 independent of the uncertainty allows for stable goal pursuit in varying levels of environmen-
 423 tal noise.

424 While many studies have explored how PFC state representations change depending on
 425 the context^{3,8,7}, there has been a lack of research focusing on strategic processes including
 426 multiple contexts and stages. We targeted to fill this gap, and our findings demonstrated
 427 that uncertainty, which can influence behavioral strategy selection^{16,17,19,13}, and goal, which
 428 affects action selection, are independently represented in the LPFC. Thus, the LPFC is
 429 capable of guiding behavioral strategies while setting and pursuing goals independently of
 430 the strategy. This suggests that the LPFC possesses the ability to establish a stable hierarchy
 431 of strategy selection and action selection.

432 In our task, we observed that when specific goals were not provided, the decrease in value
 433 difference led to a decrease in the optimality and consistency of behavior (Fig. 1e), consistent
 434 with previous findings in value-based decision-making research^{46,47,48}. Interestingly, during
 435 goal-directed learning, we observed that participants maintained optimal choices despite the
 436 influence of uncertainty on action value. According to our neural data analysis, this stable
 437 goal-pursuit behavior is supported by the neural goal robustness to uncertainty in the LPFC
 438 and OFC. Thus, our study contributes to understanding the role of the PFC in guiding
 439 multistage decision-making involving multiple contexts at the neural representational level,
 440 which was previously difficult to explain through value learning alone^{85,86,87,88}.

441 We confirmed that participants exhibiting flexible behavior also show high behavioral
 442 stability. It contrasts with previous studies widely discussing the stability-flexibility trade-
 443 off^{30,31,32,33,34,35}. However, recent task-switch-related research suggests that these two charac-

⁴⁴⁴ teristics can be controlled by independent mechanisms and are not necessarily conflicting^{89,90}.

⁴⁴⁵ The neural evidence we present could serve as a key to explaining the brain's ability to be

⁴⁴⁶ flexible yet stable, along with computational modeling research.

⁴⁴⁷ Generally, the OFC is known to compress and judge external information such as task-

⁴⁴⁸ relevant state space, confidence, emotion, and value estimation^{21,22,23,91,24,26}. Our results

⁴⁴⁹ showing that the OFC is most sensitive to uncertainty are consistent with the existing

⁴⁵⁰ literature. Furthermore, regarding specific goals as explicit task-specific information, the

⁴⁵¹ OFC seems to extract latent information about environmental uncertainty while minimizing

⁴⁵² interference between them. Thus, by factorizing abstract latent information and task-specific

⁴⁵³ information, the OFC facilitates the generalization of important information as tasks change.

⁴⁵⁴ On the other hand, the LPFC extracts important context information related to meta-

⁴⁵⁵ control, task switching, and planning^{92,17,18,93,51}, guiding strategy selection and action se-

⁴⁵⁶ lection. The observation that the LPFC is most sensitive to goal information is consistent

⁴⁵⁷ with the existing literature. Furthermore, the ability of the LPFC to independently repre-

⁴⁵⁸ sent task-relevant goals in uncertain environments enables consistent pursuit of task goals

⁴⁵⁹ across various environments. That is, the factorized representation in the LPFC enables the

⁴⁶⁰ transfer of goals across environments.

⁴⁶¹ A promising direction for future studies is to develop neural network models that integrate

⁴⁶² factorized representations of goal and uncertainty, extending approaches such as β -VAE

⁴⁶³ from sensory perception^{94,95} to context-dependent reinforcement learning. Testing whether

⁴⁶⁴ this architecture enables artificial agents to balance flexibility and stability, analogous to

⁴⁶⁵ human performance, would clarify the representational mechanisms underlying adaptive,

⁴⁶⁶ goal-directed behavior. Further investigation should probe how goal and uncertainty signals

⁴⁶⁷ evolve across extended timescales or in more complex tasks, assessing whether factorized

⁴⁶⁸ geometry in PFC persists in diverse scenarios.

⁴⁶⁹ In conclusion, our findings suggest a representational solution for achieving flexible yet

⁴⁷⁰ stable goal pursuit under uncertainty. By maintaining separate codes for goal and uncer-
⁴⁷¹ tainty, the PFC preserves goal-directed action plans while selectively adjusting behavior in
⁴⁷² response to environmental variability. This perspective extends conventional computational
⁴⁷³ accounts by emphasizing how the geometry of neural representations can link specialized
⁴⁷⁴ computations with more generalizable cognitive control. This representational account pro-
⁴⁷⁵ vides a mechanistic framework for understanding how the brain maintains goals across vary-
⁴⁷⁶ ing contexts, offering testable hypotheses for future computational and empirical research
⁴⁷⁷ on robust context-dependent learning.

478 **Methods**

479 **Participants**

480 We used the same participant dataset as the previous study¹⁷. Twenty-two subjects (all
481 right-handed, six females, mean age: 28 years, age range: 19 to 40 years) participated in
482 the experiment, and none of them had a history of neurological or psychiatric diseases. All
483 subjects gave informed consent, and the study was approved by the Institutional Review
484 Board of the California Institute of Technology.

485 One subject was excluded from the analysis since the subject consistently chose only one
486 of the two choices in stage 1 and never experienced one of the four goal states. Another sub-
487 ject was excluded from the analysis because of the exceptionally low behavioral performance;
488 the average choice optimality of that subject was less than 0.5 in stage 2, which translated
489 into worse performance than random choice.

490 **Task**

491 We used behavior and fMRI data published in a previous study¹⁷. Twenty participants
492 performed a sequential two-choice Markov decision task. In each trial, they began in a
493 common start state and made two sequential choices (by pressing left or right within 4 s) to
494 obtain a monetary reward in the form of a colored coin (red, yellow, or blue) at the end of the
495 sequence. If participants did not respond within 4 s, the computer selected a random choice
496 for them, and that trial was designated as a penalizing trial. The reward values (USD 0.40,
497 0.20, and 0.10) were randomly assigned to each coin color for each subject at the beginning
498 of the experiment.

499 Before the main experiment, each participant completed a pretraining session consisting
500 of 100 trials in which the state-transition probability was fixed at (0.5, 0.5). During these
501 trials, a white “collection box” was presented, indicating that any colored coin would yield

502 its assigned monetary reward. This pretraining was intended to allow participants ample
 503 opportunity—based on pilot data indicating that 80 trials are sufficient—to become familiar
 504 with both the sequential choice structure and the general reward contingency of a two-choice
 505 Markov task.

506 Following pretraining, the experiment proceeded in five separate scanning sessions of
 507 approximately 80 trials each, for a total of 400 trials in the main task. Each scanning
 508 session featured two conditions that manipulated the goal or collection box presented at the
 509 start of each trial. In the specific goal condition, the collection box was rendered in a single
 510 color (red, yellow, blue, or gray), indicating that only one particular coin color would be
 511 valuable on that trial (i.e., yield money if obtained). In the non-specific goal condition, the
 512 collection box was white, indicating that any of the three colored end states would provide
 513 its associated monetary outcome.

514 Throughout the main task, participants were not informed of the numeric state-transition
 515 probabilities, only that these contingencies could change. Specifically, the transitions alter-
 516 nated between (0.9, 0.1) and (0.5, 0.5) across short blocks to induce shifts in task predictabil-
 517 ity. Each block contained three to five trials under the (0.9, 0.1) condition and five to seven
 518 trials under the (0.5, 0.5) condition. This design ensured that participants experienced
 519 periods of relatively deterministic transitions and more uncertain transitions, encouraging
 520 the engagement of both model-based and model-free learning strategies. The time between
 521 states (and between trials) was sampled from a uniform distribution (1–4 s), and the reward
 522 outcome was displayed for 2 s at the end of each trial.

523 Participants were instructed that they would receive the cumulative monetary earnings
 524 from the task and that they should learn, through experience, which choices were more
 525 likely to lead to each colored coin. They were also aware that goal states and transition
 526 probabilities could vary. No further explicit information regarding probabilities or block
 527 lengths was provided, ensuring that they relied on ongoing experience to guide their choices.

528 **Behavioral measures**

529 We employed three behavioral measures to characterize goal-directed learning behavior.
 530 First, to assess behavioral flexibility between goals, we used *choice versatility*, defined as
 531 the switch in choice upon a goal change. In trials where the goal changed, the choice ver-
 532 satility was assigned a value of one if the current choice differed from the previous one at
 533 the same state, and zero if the choice remained consistent. Second, to evaluate behavioral
 534 stability within a goal, we used the *choice consistency* measure. For trials with the same
 535 goal, the choice consistency was one if the current choice was identical to the previous one at
 536 the same state, and zero if the choice changed. Lastly, *choice optimality* was used to assess
 537 the behavioral performance of goal-directed learning. For each trial, the choice optimality
 538 was assigned a value of one if the agent made an optimal value-based decision, and zero if
 539 the choice was not optimal. If two choices had identical action values, resulting in the same
 540 expected reward regardless of the choice, those trials were excluded from the analysis and
 541 choice optimality was not calculated for them.

542 **Simulation**

543 To simulate the learning processes, we generated virtual episodes using both the MB and MF
 544 learning agents. The MB agent employed both FORWARD learning⁹⁶ and BACKWARD
 545 planning¹⁷. The FORWARD learning mechanism enables the agent to learn the model
 546 of the environment by computing the state prediction error to update state-action-state
 547 transition probabilities and corresponding state-action values. BACKWARD planning allows
 548 instantaneous updates to the state-action value signal in response to changes in the goal,
 549 which defines the rewards of the outcome states. Whenever a goal is given, the model-based
 550 agent calculates the action value (Q_{MB}) using BACKWARD planning as follows. In doing so,
 551 it uses the estimated state-transition matrix $T(s, a, s')$ for the given environment to compute

552 the action value:

$$r(s) = \begin{cases} R, & \text{for a goal state,} \\ 0, & \text{otherwise.} \end{cases}$$

for $i = 3, 2$

for $s \in S_{i-1}$

$$Q_{\text{MB}}(s, a) = \sum_{s'} T(s, a, s') [r(s') + \gamma \max_{a'} Q_{\text{MB}}(s', a')], \quad \forall a$$

end for

end for

553 Here, R is the reward value corresponding to the goal state, and S_i refers to the set of
 554 states in the i -th stage. s, s' refers to the current and the next state, respectively. a, a' refers
 555 to the action in the current state and in the next state, respectively. Since γ is the temporal
 556 discount factor and, in our task, the actual reward is only given at the final stage, we set
 557 $\gamma = 1$.

558 In addition, at each state transition, it calculates a state prediction error (SPE) to update
 559 $T(s, a, s')$:

$$\delta_{\text{SPE}} = 1 - T(s, a, s'),$$

$$\Delta T(s, a, s') = \eta \delta_{\text{SPE}},$$

$$Q_{\text{MB}}(s, a) = \sum_{s'} T(s, a, s') [r(s') + \gamma \max_{a'} Q_{\text{MB}}(s', a')],$$

560 where η is the learning rate of the state-transition probability estimation.

561 The MF learning agent was implemented using SARSA, which utilized conventional

⁵⁶² temporal-difference (TD) updates to compute the reward prediction error for state-action
⁵⁶³ value (Q_{MF}) updates⁴⁵:

$$\delta_{\text{RPE}} = r(s') + \gamma Q_{\text{MF}}(s', a') - Q_{\text{MF}}(s, a),$$

$$\Delta Q_{\text{MF}}(s, a) = \alpha \delta_{\text{RPE}},$$

⁵⁶⁴ where α is the learning rate of Q_{MF} .

⁵⁶⁵ Both the MB and MF models then compute the choice probabilities from action values
⁵⁶⁶ using the softmax function, and their stochastic choices constitute the simulation behavior:

$$P(a | s) = \frac{\exp(\tau Q(s, a))}{\sum_{a'} \exp(\tau Q(s, a'))}.$$

⁵⁶⁷ Here, τ serves as the inverse temperature controlling how strongly the model exploits
⁵⁶⁸ value differences.

⁵⁶⁹ To maintain consistency with human behavioral data, we utilized the total number of
⁵⁷⁰ trials, block condition sequences, and specific goal sequences directly from the data of 20
⁵⁷¹ human participants. The RL agents made their own choices, leading to state transitions
⁵⁷² that differed from those in the human data. For each human experimental sequence, we
⁵⁷³ performed simulations of both MB and MF agents using 1000 different random parameter
⁵⁷⁴ sets, resulting in a total of 20,000 samples for analysis. Both the MB and MF models include
⁵⁷⁵ two free parameters—a learning rate and a softmax inverse temperature. Specifically, the
⁵⁷⁶ MB model uses the SPE learning rate η , whereas the MF model uses the RPE learning rate
⁵⁷⁷ α .

578 **fMRI data collection and pre-processing**

579 We used the fMRI dataset provided by the previous study¹⁷. MRI images were obtained
 580 from the Caltech Brain Imaging Center, which uses a 3T Siemens (Erlangen) Trio scanner
 581 with a 32-channel radiofrequency coil. Structural images were collected using a standard
 582 MPRAGE pulse sequence (long repetition time (TR): 1,500 ms, short echo time (TE): 2.63
 583 ms, flip angle: 10°, voxel size: 1 mm × 1 mm × 1 mm). For the functional images, 45 slices
 584 were collected at an angle of 30° from the anterior commissure-posterior commissure axis
 585 using a one-shot echo-planar imaging pulse sequence (TR: 2,800 ms, TE: 30 ms, flip angle:
 586 80°, field of view: 100 mm, voxel size: 3 mm × 3 mm × 3 mm).

587 fMRI data were preprocessed using the SPM8 software package. Preprocessing steps
 588 were conducted for each participant individually. Slice-timing correction was applied to
 589 adjust for acquisition time differences across slices within each image, using the first slice
 590 as the reference. To correct for participant motion, realignment was performed with the
 591 mean of the images as the reference. Each participant's structural image was coregistered to
 592 the mean functional realigned image and normalized to the Montreal Neurological Institute
 593 (MNI) 152 template. The functional images were subsequently spatially transformed based
 594 on these normalization parameters, aligning them to the MNI152 template brain to account
 595 for anatomical variability across participants.

596 For additional preprocessing steps required for multivoxel pattern analysis (MVPA), we
 597 utilized the Princeton MVPA toolbox (<http://code.google.com/p/princeton-mvpa-toolbox>)
 598 and custom code. Within each scanning session, the fMRI data were detrended, and the
 599 BOLD time series of each voxel was z-scored. The resulting multi-voxel time series were
 600 then used as trial-by-trial brain activity patterns.

601 We performed ROI analysis based on eight brain regions: the ventrolateral prefrontal cor-
 602 tex (vlPFC), dorsolateral prefrontal cortex (dlPFC), orbitofrontal cortex (OFC), anterior cin-
 603 gulate cortex (ACC), pre-supplementary motor area (preSMA), primary visual cortex (V1),

604 hippocampus (HPC), and ventral striatum (vStr) (Supplementary Fig. 2). All fMRI data
 605 used in the analyses were extracted from brain regions defined by the automated anatomical
 606 labeling (AAL3) atlas⁶², except for the preSMA. The preSMA was defined as the preSMA
 607 region from the JuBrain Anatomy toolbox⁹⁷ (the SPM Anatomy Toolbox; <https://www.fz-juelich.de/en/inm/inm-7/resources/tools/jubrain-anatomy-toolbox>). The vIPFC was de-
 608 fined as the triangular part of the inferior frontal gyrus¹⁷. The dlPFC was defined as the
 609 middle frontal gyrus. The OFC was defined as bilateral inferior, middle, and superior orbital
 610 gyri and bilateral rectal gyri²³. The ACC was defined as the pregenual and supracallosal an-
 611 terior cingulate cortex. The V1 was defined as calcarine fissures and the surrounding cortex.
 612 The hippocampus and ventral striatum were defined as the AAL3 ROIs with the same name,
 613 respectively. In our preliminary analyses, we used under-sampling to match the number of
 614 voxels between ROIs and confirmed that differences in voxel number do not significantly
 615 affect our results. All results were bilaterally averaged for each region since there were no
 616 significant differences between hemispheres.

618 fMRI decoding analyses

619 For all our analyses (simple decoding, shattering, and CCGP), we trained linear SVMs
 620 on multivoxel patterns from a participant's ROI to separate task-variable classes (for simple
 621 decoding analysis; Fig. 2) or particular dichotomies (for shattering analysis; Fig. 3 and CCGP
 622 analysis; Fig. 4). We treated the pre-processed voxel-wise BOLD time courses within each
 623 ROI as trial-by-trial neural activity patterns. To avoid session-related dependencies between
 624 the training and test data, we used leave-one-session-out cross-validation for all analyses.
 625 For each participant, the BOLD patterns from one of the five fMRI scanning sessions served
 626 as the test set, whereas the data from the remaining sessions formed the training set used
 627 to fit the classifier. All reported classification results correspond to the mean test accuracy
 628 obtained across the five cross-validation folds. Additionally, the label imbalance effect was

629 ruled out by under-sampling of the larger label. To minimize the random effect of the
630 under-sampling, we repeated sampling 100 times with different random seeds, and the test
631 accuracies of all 100 samplings were averaged for the main analysis.

632 To measure trial-event-specific neural measures, we used fMRI data scanned in the cor-
633 responding time bin. We recorded the timings of the following trial events: Fixation 1,
634 Stimulus 1 (representing the first state S1), Choice 1 (A1), Fixation 2, Stimulus 2 (S2),
635 Choice 2 (A2), Fixation 3, and Stimulus 3 (S3). The timings of the eight events were
636 recorded for each trial. Regarding the sluggish hemodynamic response, we used the first
637 fMRI volume recorded immediately after the occurrence of a specific event as the neural
638 response elicited by that event. However, choices 1 and 2 were always followed by instant
639 fixation cues (Fixations 2 and 3, respectively). As a result, each fixation cue after a choice
640 always preceded the scanning of the choice-specific response volume. Therefore, we labeled
641 volumes scanned after Fixation 2 and 3 as choice-specific activity (A1 and A2, respectively).

642 We performed decoding analyses on voxel-level multivariate patterns from ROIs of each
643 participant. Decoding accuracy was calculated as a single value per participant based on
644 multiple trials. The group-level statistical tests were performed, treating each participant
645 as a random sample. No multiple comparisons correction was performed since we tested
646 pre-defined ROIs' results individually and did not perform an exploratory search for some
647 specific effect using multiple samples.

648 The association between a neural measure and a behavioral measure was assessed using
649 Pearson's correlation coefficient across participants. Here, investigating the relationship
650 between neural and behavioral measures via correlation were exploratory. Therefore, we
651 performed correlation analyses across all eight ROIs and corrected for multiple comparisons
652 using the Benjamini-Hochberg procedure (FDR, $q = 0.05$). All correlation results (Fig. 3c,
653 and 4b) reflect this correction, and the adjusted p-values are reported.

654 **Shattering analysis**

655 To investigate the representational geometry of goal and uncertainty in multi-voxel patterns
 656 of brain regions, we computed the shattering dimensionality (SD)^{11,41} by averaging test
 657 accuracies of all linear support vector machines (SVM) trained to dichotomize task variables.
 658 Accordingly, the SD quantifies the separability of neural embeddings associated with each
 659 class of a task variable. The number of total dichotomies, which is equivalent to the number
 660 of ways of binary labeling, is determined by the number of task variable classes. There are
 661 2^C ways of binary labeling with the C classes. The actual number of dichotomies required to
 662 be tested reduces to $2^{C-1} - 1$ after excepting the two cases of all positive or negative labeling
 663 and removing half of the duplicated cases due to the symmetry of binary labeling. Since it
 664 is based on binary classifiers, the chance level is 0.5.

665 To investigate the mixed embedding structure of specific goals and uncertainty, we per-
 666 formed a shattering analysis on goal-uncertainty combined classes (red-low, blue-low, yellow-
 667 low, red-high, blue-high, yellow-high). We categorized all dichotomies based on six classes
 668 into four types of classification: goal, uncertainty, linear, and nonlinear (Supplementary Fig.
 669 6). The goal type included three dichotomies that separated each goal class from the others
 670 (e.g., red-low & red-high vs. the other classes), and the goal SD was defined as the average
 671 test accuracy of these three dichotomies (Fig. 3b). Similarly, the uncertainty type included
 672 one dichotomy separating the low vs. high uncertainty conditions (3 classes vs. 3 classes),
 673 and its test accuracy defined the uncertainty SD.

To determine linear and nonlinear type dichotomies, we assessed linear separability in a random dataset with a linearly mixed representation. Neural activity due to goal and uncertainty was expressed as:

$$\mathbf{y} = W_g \mathbf{g} + W_u \mathbf{u} + \mathbf{b}$$

674 where \mathbf{y} is the N-dimensional neural response, \mathbf{g} is a 3×1 one-hot vector representing
 675 one of the three specific goals, u is a binary scalar variable representing uncertainty, W_g
 676 and W_u are the linear weight parameters for the two variables, and \mathbf{b} is the N-dimensional
 677 bias parameter independent of goal and uncertainty. Using $N=3$, we generated 1000 sets
 678 of random parameters and assessed the linear separability of all dichotomies for the six
 679 classes. MATLAB's 'perceptron' function was used for linear classification. Nine dichotomies
 680 not linearly separable across all random seeds were categorized into the nonlinear type.
 681 Dichotomies not categorized into the goal, uncertainty, or nonlinear types were categorized
 682 into the linear type.

683 Cross-condition generalization performance

684 To examine how the context condition (uncertainty) influences the neural embedding of
 685 the goal, we quantified cross-condition generalization performance (CCGP), defined as the
 686 generalized accuracy of a linear decoder across contexts. Because the uncertainty condition is
 687 binary, two cross-decoding directions arise: (i) training on low-uncertainty trials and testing
 688 on high-uncertainty trials, and (ii) the reverse. The CCGP is defined as the mean test
 689 accuracy obtained from these two directions.

690 A high CCGP indicates that the decoded variable (specific goals) is disentangled from
 691 contextual condition (uncertainty levels). If the goal representation is modulated by uncer-
 692 tainty, the CCGP will be lower than the within-condition decoding accuracy. To capture
 693 the performance decrease by comparing CCGP against the SD, we performed the cross-
 694 uncertainty goal decoding analysis consistent with the shattering analysis. Specifically, we
 695 implemented three linear dichotomies that discriminate each specific goal from the other two
 696 (red vs. blue, yellow; blue vs. red, yellow; yellow vs. red, blue) and averaged their test
 697 accuracies.

698 Consistent with all other classification analyses, we adopted a leave-one-session-out val-

699 idation scheme: each fMRI scanning session served as the test set, while the remaining
 700 sessions constituted the training set. For CCGP, this procedure was repeated in both cross-
 701 condition directions (e.g., training on low-uncertainty data from the remaining sessions and
 702 testing on high-uncertainty data from the held-out session, and vice versa). The CCGP
 703 value is the mean of the two resulting test accuracies.

704 All other analysis settings matched those used in the shattering-dimension analysis.

705 Statistical analysis

706 Unless otherwise stated, hypothesis-driven tests were evaluated at $\alpha = 0.05$ (two-tailed). Be-
 707 cause each ROI was selected *a priori* on theoretical grounds, statistical tests were performed
 708 independently for every ROI without a family-wise correction across ROIs (Fig. 2, Fig. 3b,
 709 Fig. 4a)^{73,98,99}. Where a single test involved several task events or SD categories within
 710 the *same* ROI, these comparisons were likewise planned and reported without additional
 711 correction.

712 In contrast, the correlations between neural metrics and behavior (Fig. 3c–d, Fig. 4b)
 713 were exploratory. Here we corrected across the eight ROIs using the Benjamini–Hochberg
 714 false-discovery-rate procedure ($q = 0.05$) and report the adjusted q -values. Exact p - and
 715 q -values are provided in Supplementary Tables 2–4.

⁷¹⁶ **Data availability**

⁷¹⁷ The human behavioral data used in this study are available in the GitHub repository at
⁷¹⁸ <https://github.com/brain-machine-intelligence/RLdim-mvpa-model>. The processed fMRI
⁷¹⁹ data generated in this study (ROI-masked EPI) have been deposited in the OSF database
⁷²⁰ (<https://osf.io/2gyue>). The data used to create the figures in this paper are provided in the
⁷²¹ Source Data file. Source data are provided with this paper.

⁷²² **Code availability**

⁷²³ The code used for the neural and simulation analyses in this study is available in the GitHub
⁷²⁴ repository (<https://github.com/brain-machine-intelligence/RLdim-mvpa-model>) and archived
⁷²⁵ at Zenodo (<https://doi.org/10.5281/zenodo.17412741>)¹⁰⁰.

726 **References**

727 **References**

- 728 1. Miller, E. K. & Cohen, J. D. An Integrative Theory of Prefrontal Cortex Function.
729 *Annual Review of Neuroscience* **24**, 167–202 (2001).
- 730 2. Cohen, J. D. Cognitive Control. In *The Wiley Handbook of Cognitive Control*, chap. 1,
731 1–28 (John Wiley & Sons, Ltd, 2017).
- 732 3. Mante, V., Sussillo, D., Shenoy, K. V. & Newsome, W. T. Context-dependent compu-
733 tation by recurrent dynamics in prefrontal cortex. *Nature* **503**, 78–84 (2013).
- 734 4. Stokes, M. G. *et al.* Dynamic Coding for Cognitive Control in Prefrontal Cortex. *Neuron*
735 **78**, 364–375 (2013).
- 736 5. Cromer, J. A., Roy, J. E. & Miller, E. K. Representation of Multiple, Independent
737 Categories in the Primate Prefrontal Cortex. *Neuron* **66**, 796–807 (2010).
- 738 6. Aoi, M. C., Mante, V. & Pillow, J. W. Prefrontal cortex exhibits multidimensional
739 dynamic encoding during decision-making. *Nature Neuroscience* **23**, 1410–1420 (2020).
- 740 7. Flesch, T., Juechems, K., Dumbalska, T., Saxe, A. & Summerfield, C. Orthogonal
741 representations for robust context-dependent task performance in brains and neural
742 networks. *Neuron* **110**, 1258–1270.e11 (2022).
- 743 8. Castegnetti, G., Zurita, M. & De Martino, B. How usefulness shapes neural represen-
744 tations during goal-directed behavior. *Science Advances* **7**, eabd5363 (2021).
- 745 9. Molinaro, G. & Collins, A. G. E. A goal-centric outlook on learning. *Trends in Cognitive
746 Sciences* (2023).

- 747 10. Muhle-Karbe, P. S. *et al.* Goal-seeking compresses neural codes for space in the human
748 hippocampus and orbitofrontal cortex. *Neuron* **111**, 3885–3899.e6 (2023).
- 749 11. Rigotti, M. *et al.* The importance of mixed selectivity in complex cognitive tasks.
750 *Nature* **497**, 585–590 (2013).
- 751 12. Huettel, S. A., Song, A. W. & McCarthy, G. Decisions under Uncertainty: Probabilistic
752 Context Influences Activation of Prefrontal and Parietal Cortices. *Journal of Neuroscience* **25**, 3304–3311 (2005).
- 753 13. Soltani, A. & Izquierdo, A. Adaptive learning under expected and unexpected uncertainty.
754 *Nature Reviews Neuroscience* **20**, 635–644 (2019).
- 755 14. Hsu, M., Bhatt, M., Adolphs, R., Tranel, D. & Camerer, C. F. Neural Systems Responding to Degrees of Uncertainty in Human Decision-Making. *Science* **310**, 1680–
756 1683 (2005).
- 757 15. Soltani, A. & Koechlin, E. Computational models of adaptive behavior and prefrontal
758 cortex. *Neuropsychopharmacology* **47**, 58–71 (2022).
- 759 16. Daw, N. D., Niv, Y. & Dayan, P. Uncertainty-based competition between prefrontal and
760 dorsolateral striatal systems for behavioral control. *Nature Neuroscience* **8**, 1704–1711
761 (2005).
- 762 17. Lee, S. W., Shimojo, S. & O'Doherty, J. P. Neural Computations Underlying Arbitration
763 between Model-Based and Model-free Learning. *Neuron* **81**, 687–699 (2014).
- 764 18. Kim, D., Jeong, J. & Lee, S. W. Prefrontal solution to the bias-variance tradeoff during
765 reinforcement learning. *Cell Reports* **37**, 110185 (2021).
- 766 19. Collins, A. G. E. & Cockburn, J. Beyond dichotomies in reinforcement learning. *Nature
767 Reviews Neuroscience* **21**, 576–586 (2020).

- 770 20. Smittenaar, P., FitzGerald, T. H. B., Romei, V., Wright, N. D. & Dolan, R. J. Disrup-
771 tion of Dorsolateral Prefrontal Cortex Decreases Model-Based in Favor of Model-free
772 Control in Humans. *Neuron* **80**, 914–919 (2013).
- 773 21. De Martino, B., Fleming, S. M., Garrett, N. & Dolan, R. J. Confidence in value-based
774 choice. *Nature Neuroscience* **16**, 105–110 (2013).
- 775 22. Lebreton, M., Abitbol, R., Daunizeau, J. & Pessiglione, M. Automatic integration of
776 confidence in the brain valuation signal. *Nature Neuroscience* **18**, 1159–1167 (2015).
- 777 23. Schuck, N. W., Cai, M. B., Wilson, R. C. & Niv, Y. Human Orbitofrontal Cortex
778 Represents a Cognitive Map of State Space. *Neuron* **91**, 1402–1412 (2016).
- 779 24. Park, S. A., Miller, D. S., Nili, H., Ranganath, C. & Boorman, E. D. Map Making:
780 Constructing, Combining, and Inferring on Abstract Cognitive Maps. *Neuron* **107**,
781 1226–1238.e8 (2020).
- 782 25. Whittington, J. C. R., McCaffary, D., Bakermans, J. J. W. & Behrens, T. E. J. How
783 to build a cognitive map. *Nature Neuroscience* **25**, 1257–1272 (2022).
- 784 26. Padoa-Schioppa, C. & Assad, J. A. Neurons in the orbitofrontal cortex encode economic
785 value. *Nature* **441**, 223–226 (2006).
- 786 27. O'doherty, J. P. Lights, Camembert, Action! The Role of Human Orbitofrontal Cortex
787 in Encoding Stimuli, Rewards, and Choices. *Annals of the New York Academy of
788 Sciences* **1121**, 254–272 (2007).
- 789 28. Rigotti, M., Rubin, D. B. D., Morrison, S., Salzman, C. & Fusi, S. Attractor concretion
790 as a mechanism for the formation of context representations. *Neuroimage* **52**, 833–847
791 (2010).

- 792 29. Saez, A., Rigotti, M., Ostojic, S., Fusi, S. & Salzman, C. Abstract context representa-
793 tions in primate amygdala and prefrontal cortex. *Neuron* **87**, 869–881 (2015).
- 794 30. Goschke, T. Volition in Action: Intentions, Control Dilemmas, and the Dynamic Reg-
795 ulation of Cognitive Control (2013).
- 796 31. Hommel, B. Chapter Two - Between Persistence and Flexibility: The Yin and Yang
797 of Action Control. In Elliot, A. J. (ed.) *Advances in Motivation Science*, vol. 2, 33–67
798 (Elsevier, 2015).
- 799 32. Armbruster-Genç, D. J. N., Ueltzhöffer, K. & Fiebach, C. J. Brain Signal Variability
800 Differentially Affects Cognitive Flexibility and Cognitive Stability. *Journal of Neuro-
801 science* **36**, 3978–3987 (2016).
- 802 33. Dreisbach, G. & Fröber, K. On How to Be Flexible (or Not): Modulation of the
803 Stability-Flexibility Balance. *Current Directions in Psychological Science* **28**, 3–9
804 (2019).
- 805 34. Musslick, S. & Cohen, J. D. Rationalizing constraints on the capacity for cognitive
806 control. *Trends in Cognitive Sciences* **25**, 757–775 (2021).
- 807 35. Qiao, L., Zhang, L. & Chen, A. Control dilemma: Evidence of the stability–flexibility
808 trade-off. *International Journal of Psychophysiology* **191**, 29–41 (2023).
- 809 36. Raposo, D., Kaufman, M. T. & Churchland, A. K. A category-free neural population
810 supports evolving demands during decision-making. *Nature Neuroscience* **17**, 1784–
811 1792 (2014).
- 812 37. Fusi, S., Miller, E. K. & Rigotti, M. Why neurons mix: High dimensionality for higher
813 cognition. *Current Opinion in Neurobiology* **37**, 66–74 (2016).

- 814 38. Tang, E. *et al.* Effective learning is accompanied by high-dimensional and efficient
815 representations of neural activity. *Nature Neuroscience* **22**, 1000–1009 (2019).
- 816 39. Sheng, J. *et al.* Higher-dimensional neural representations predict better episodic mem-
817 ory. *Science Advances* **8**, eabm3829 (2022).
- 818 40. Mack, M. L., Preston, A. R. & Love, B. C. Ventromedial prefrontal cortex compression
819 during concept learning. *Nature Communications* **11**, 46 (2020).
- 820 41. Bernardi, S. *et al.* The Geometry of Abstraction in the Hippocampus and Prefrontal
821 Cortex. *Cell* **183**, 954–967.e21 (2020).
- 822 42. Chung, S. & Abbott, L. F. Neural population geometry: An approach for understanding
823 biological and artificial neural networks. *Current Opinion in Neurobiology* **70**, 137–144
824 (2021).
- 825 43. Badre, D., Bhandari, A., Keglovits, H. & Kikumoto, A. The dimensionality of neural
826 representations for control. *Current Opinion in Behavioral Sciences* **38**, 20–28 (2021).
- 827 44. Jazayeri, M. & Ostoicic, S. Interpreting neural computations by examining intrinsic
828 and embedding dimensionality of neural activity. *Current Opinion in Neurobiology* **70**,
829 113–120 (2021).
- 830 45. Sutton, R. S. & Barto, A. G. *Reinforcement Learning, Second Edition: An Introduction*
831 (MIT Press, 2018).
- 832 46. Solway, A. & Botvinick, M. M. Evidence integration in model-based tree search. *Pro-
833 ceedings of the National Academy of Sciences* **112**, 11708–11713 (2015).
- 834 47. Krajbich, I. & Rangel, A. Multialternative drift-diffusion model predicts the relation-
835 ship between visual fixations and choice in value-based decisions. *Proceedings of the
836 National Academy of Sciences* **108**, 13852–13857 (2011).

- 837 48. Hare, T. A., Schultz, W., Camerer, C. F., O'Doherty, J. P. & Rangel, A. Transformation
838 of stimulus value signals into motor commands during simple choice. *Proceedings of the*
839 *National Academy of Sciences* **108**, 18120–18125 (2011).
- 840 49. Kim, D., Park, G. Y., O'Doherty, J. P. & Lee, S. W. Task complexity interacts with
841 state-space uncertainty in the arbitration between model-based and model-free learning.
842 *Nature Communications* **10**, 5738 (2019).
- 843 50. Heo, S., Sung, Y. & Lee, S. W. Effects of subclinical depression on prefrontal–striatal
844 model-based and model-free learning. *PLOS Computational Biology* **17**, e1009003
845 (2021. 5. 14.).
- 846 51. Mushiake, H., Saito, N., Sakamoto, K., Itoyama, Y. & Tanji, J. Activity in the Lateral
847 Prefrontal Cortex Reflects Multiple Steps of Future Events in Action Plans. *Neuron*
848 **50**, 631–641 (2006).
- 849 52. Yamagata, T., Nakayama, Y., Tanji, J. & Hoshi, E. Distinct Information Representa-
850 tion and Processing for Goal-Directed Behavior in the Dorsolateral and Ventrolateral
851 Prefrontal Cortex and the Dorsal Premotor Cortex. *Journal of Neuroscience* **32**, 12934–
852 12949 (2012).
- 853 53. Kerns, J. G. *et al.* Anterior cingulate conflict monitoring and adjustments in control.
854 *Science* **303**, 1023–1026 (2004).
- 855 54. Behrens, T. E. J., Woolrich, M. W., Walton, M. E. & Rushworth, M. F. S. Learning
856 the value of information in an uncertain world. *Nature Neuroscience* **10**, 1214–1221
857 (2007).
- 858 55. Nachev, P., Wydell, H., O'Neill, K., Husain, M. & Kennard, C. The role of the pre-
859 supplementary motor area in the control of action. *NeuroImage* **36**, T155–T163 (2007).

- 860 56. O'Keefe, J. & Nadel, L. *The Hippocampus as a Cognitive Map* (Clarendon Press, 1978).
- 861 57. Crivelli-Decker, J. *et al.* Goal-oriented representations in the human hippocampus
862 during planning and navigation. *Nature Communications* **14**, 2946 (2023).
- 863 58. Stachenfeld, K. L., Botvinick, M. M. & Gershman, S. J. The hippocampus as a pre-
864 dictive map. *Nature Neuroscience* **20**, 1643–1653 (2017).
- 865 59. Pessiglione, M., Seymour, B., Flandin, G., Dolan, R. J. & Frith, C. D. Dopamine-
866 dependent prediction errors underpin reward-seeking behaviour in humans. *Nature*
867 **442**, 1042–1045 (2006).
- 868 60. Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P. & Dolan, R. J. Model-Based
869 Influences on Humans' Choices and Striatal Prediction Errors. *Neuron* **69**, 1204–1215
870 (2011).
- 871 61. Payzan-LeNestour, E., Dunne, S., Bossaerts, P. & O'Doherty, J. P. The Neural Rep-
872 presentation of Unexpected Uncertainty during Value-Based Decision Making. *Neuron*
873 **79**, 191–201 (2013).
- 874 62. Rolls, E. T., Huang, C.-C., Lin, C.-P., Feng, J. & Joliot, M. Automated anatomical
875 labelling atlas 3. *NeuroImage* **206**, 116189 (2020).
- 876 63. Ito, T. *et al.* Compositional generalization through abstract representations in human
877 and artificial neural networks. *Advances in Neural Information Processing Systems* **35**,
878 32225–32239 (2022).
- 879 64. Vapnik, V. N. & Chervonenkis, A. Ya. On the Uniform Convergence of Relative Fre-
880 quencies of Events to Their Probabilities. In Vovk, V., Papadopoulos, H. & Gammer-
881 man, A. (eds.) *Measures of Complexity: Festschrift for Alexey Chervonenkis*, 11–30
882 (Springer International Publishing, Cham, 2015).

- 883 65. Abu-Mostafa, Y. S. The Vapnik-Chervonenkis Dimension: Information versus Com-
plexity in Learning. *Neural Computation* **1**, 312–317 (1989).
- 884
885 66. Feher Da Silva, C., Lombardi, G., Edelson, M. & Hare, T. A. Rethinking model-based
and model-free influences on mental effort and striatal prediction errors. *Nature Human
Behaviour* (2023).
- 886
887
888 67. Möhring, L. & Gläscher, J. Prediction errors drive dynamic changes in neural patterns
that guide behavior. *Cell Reports* **42** (2023). URL [https://www.cell.com/cell-reports/abstract/S2211-1247\(23\)00942-7](https://www.cell.com/cell-reports/abstract/S2211-1247(23)00942-7).
- 889
890
891 68. Kool, W., Gershman, S. J. & Cushman, F. A. Cost-benefit arbitration between multiple
reinforcement-learning systems. *Psychological Science* **28**, 1321–1333 (2017).
- 892
893 69. Kim, T. *et al.* Neurocomputational model of compulsion: deviating from an uncertain
goal-directed system. *Brain* **147**, 2230–2244 (2024).
- 894
895 70. Kiani, R. & Shadlen, M. N. Representation of confidence associated with a decision by
neurons in the parietal cortex. *Science* **324**, 759–764 (2009).
- 896
897 71. Hebart, M. N., Schriever, Y., Donner, T. H. & Haynes, J.-D. The relationship between
perceptual decision variables and confidence in the human brain. *Cerebral Cortex* **26**,
118–130 (2016).
- 898
899
900 72. Gherman, S. & Philiastides, M. G. Human vmpfc encodes early signatures of confidence
in perceptual decisions. *eLife* **7**, e38293 (2018).
- 901
902 73. van Bergen, R. S., Ji Ma, W., Pratte, M. S. & Jehee, J. F. M. Sensory uncertainty
decoded from visual cortex predicts behavior. *Nature Neuroscience* **18**, 1728–1730
(2015).
- 903
904

- 905 74. Li, H.-H., Sprague, T. C., Yoo, A. H., Ma, W. J. & Curtis, C. E. Joint representation of
906 working memory and uncertainty in human cortex. *Neuron* **109**, 3699–3712.e6 (2021).
- 907 75. Geurts, L. S., Cooke, J. R. H., van Bergen, R. S. & Jehee, J. F. M. Subjective confidence
908 reflects representation of bayesian probability in cortex. *Nature Human Behaviour* **6**,
909 294–305 (2022).
- 910 76. Friedman, N. P. & Robbins, T. W. The role of prefrontal cortex in cognitive control
911 and executive function. *Neuropsychopharmacology* **47**, 72–89 (2022).
- 912 77. Shenhav, A., Botvinick, M. M. & Cohen, J. D. The Expected Value of Control: An
913 Integrative Theory of Anterior Cingulate Cortex Function. *Neuron* **79**, 217–240 (2013).
- 914 78. Dorfman, H. M. & Gershman, S. J. Controllability governs the balance between pavlo-
915 vian and instrumental action selection. *Nature Communications* **10**, 5826 (2019).
- 916 79. Liljeholm, M., Dunne, S. & O'Doherty, J. P. Differentiating neural systems mediating
917 the acquisition vs. expression of goal-directed and habitual behavioral control. *European*
918 *Journal of Neuroscience* **41**, 1358–1371 (2015).
- 919 80. Baram, A. B., Muller, T. H., Nili, H., Garvert, M. M. & Behrens, T. E. J. Entorhinal
920 and ventromedial prefrontal cortices abstract and generalize the structure of reinforce-
921 ment learning problems. *Neuron* **109**, 713–723.e7 (2021).
- 922 81. Witkowski, P. P., Park, S. A. & Boorman, E. D. Neural mechanisms of credit assignment
923 for inferred relationships in a structured world. *Neuron* **110**, 2680–2690.e9 (2022).
- 924 82. Ritz, H. & Shenhav, A. Orthogonal neural encoding of targets and distractors supports
925 multivariate cognitive control. *Nature Human Behaviour* 1–17 (2024).

- 926 83. Dekker, R. B., Otto, F. & Summerfield, C. Curriculum learning for human com-
927 positional generalization. *Proceedings of the National Academy of Sciences* **119**,
928 e2205582119 (2022).
- 929 84. Flesch, T., Saxe, A. & Summerfield, C. Continual task learning in natural and artificial
930 agents. *Trends in Neurosciences* **46**, 199–210 (2023).
- 931 85. Juechems, K. & Summerfield, C. Where Does Value Come From? *Trends in Cognitive
932 Sciences* **23**, 836–850 (2019).
- 933 86. Averbeck, B. & O'Doherty, J. P. Reinforcement-learning in fronto-striatal circuits.
934 *Neuropsychopharmacology* **47**, 147–162 (2022).
- 935 87. Martino, B. D. & Cortese, A. Goals, usefulness and abstraction in value-based choice.
936 *Trends in Cognitive Sciences* **27**, 65–80 (2023).
- 937 88. Frömer, R. & Shenhav, A. Filling the gaps: Cognitive control as a critical lens for un-
938 derstanding mechanisms of value-based decision-making. *Neuroscience & Biobehavioral
939 Reviews* **134**, 104483 (2022).
- 940 89. Egner, T. Principles of cognitive control over task focus and task switching. *Nature
941 Reviews Psychology* 1–13 (2023).
- 942 90. Nack, C. & Yu-Chin, C. Cognitive flexibility and stability at the task-set level: A
943 dual-dimension framework. *advances.in/psychology* **1**, 1–28 (2023).
- 944 91. O'Reilly, R. C. Unraveling the Mysteries of Motivation. *Trends in Cognitive Sciences*
945 **24**, 425–434 (2020).
- 946 92. Domenech, P. & Koechlin, E. Executive control and decision-making in the prefrontal
947 cortex. *Current Opinion in Behavioral Sciences* **1**, 101–106 (2015).

- 948 93. Vaidya, A. R. & Badre, D. Abstract task representations for inference and control.
949 *Trends in Cognitive Sciences* **26**, 484–498 (2022).
- 950 94. Higgins, I. *et al.* beta-vae: Learning basic visual concepts with a constrained variational
951 framework (2017). URL <https://openreview.net/forum?id=Sy2fzU9gl>.
- 952 95. Higgins, I. *et al.* Unsupervised deep learning identifies semantic disentanglement in
953 single inferotemporal face patch neurons. *Nature Communications* **12**, 6456 (2021).
- 954 96. Gläscher, J., Daw, N., Dayan, P. & O'Doherty, J. P. States versus Rewards: Dissociable
955 Neural Prediction Error Signals Underlying Model-Based and Model-Free Reinforce-
956 ment Learning. *Neuron* **66**, 585–595 (2010).
- 957 97. Eickhoff, S. B. *et al.* Assignment of functional activations to probabilistic cytoarchitec-
958 tonic areas revisited. *NeuroImage* **36**, 511–521 (2007).
- 959 98. Gadassi Polack, R., Mollick, J. A., Keren, H., Joormann, J. & Watts, R. Neural
960 responses to reward valence and magnitude from pre- to early adolescence. *NeuroImage*
961 **275**, 120166 (2023).
- 962 99. Kenyon, K. H. *et al.* The characteristics and reproducibility of motor speech
963 functional neuroimaging in healthy controls. *Frontiers in Human Neuroscience*
964 **18** (2024). URL <https://www.frontiersin.org/https://www.frontiersin.org/journals/human-neuroscience/articles/10.3389/fnhum.2024.1382102/full>.
- 966 100. Sung, Y., Rigotti, M. & Lee, S. W. Factorized embedding of goal and uncertainty in
967 the lateral prefrontal cortex guides stably flexible learning (2025). URL <https://doi.org/10.5281/zenodo.17412741>.

969 **Acknowledgements**

970 We sincerely thank Jungwon Ryu, Taekwan Kim, and Takuya Ito for discussions that helped
971 refine the direction of this research. We also appreciate Yujin Cha, Jaehoon Shin, and Heejun
972 Kim for their broad and consistent support for this project. We deeply thank Sungyoung
973 Kim for meaningful comments and valuable help in finalizing this work. This research was
974 supported by the National Research Foundation of Korea (NRF), funded by the Korean
975 government (MSIT) (No. RS-2024-00439903).

976 **Author contributions**

977 Conceptualization, Y.S. and S.W.L.; Methodology, Y.S., M.R, and S.W.L.; Software. Y.S.;
978 Validation, Y.S.; Formal Analysis, Y.S.; Investigation, Y.S. and S.W.L.; Resources, S.W.L.;
979 Data Curation, S.W.L.; Writing – Original Draft, Y.S.; Writing – Review & Editing, Y.S.,
980 M.R., and S.W.L.; Visualization, Y.S.; Supervision, S.W.L.; Project Administration, S.W.L;
981 Funding Acquisition, S.W.L

982 **Competing interests**

983 The authors declare no competing interests.

984 **Figure Legends/Captions**

Fig. 1 | Flexible and stable human behavior during goal-directed learning

(a) Task structure as a binary decision tree. S_k and A_k denote the state and action at stage k respectively ($k = 1, 2, 3$). Each trial begins with the same initial state. (b) Block conditions that determine the task context. At the start of each trial, the specific goal condition is represented by one of the boxes (red, blue, yellow) alongside a fractal image indicating the state. In the non-specific goal condition, a white box is shown with the fractal image. Participants were informed that state-transition probabilities could change, but specific probabilities were not provided, and no direct cues about the uncertainty condition were given during the experiment. (c) Predicted impact of uncertainty on action value difference between left and right choices. The left side of the panel shows the value difference calculated for a specific context and state. The right side shows the average value difference across all contexts and states. (d) Effect of goal and uncertainty on the performance of human participants ($n = 20$). Points next to each boxplot represent individual human participants. Asterisks denote statistical significance (paired t-test, ***: $p < 0.0001$). For the box plots, the center lines, box limits, and whiskers represent medians, upper/lower quartiles, and $1.5 \times$ interquartile ranges, respectively. All statistical tests were two-sided. See Supplementary Table 1 for full statistical information. (e) Behavioral measures defined within the specific goal condition. Each action in a trial is assigned a value of 0 or 1 based on the criteria. (f) Scatter plot of the behavioral measures of human participants and virtual RL agents. For human results, each point represents a single participant. For simulation results, each point represents an agent, and MB and MF agents were simulated with 20,000 sets of random parameters, respectively (Simulation). See Supplementary Table 1 for the statistical details of the Pearson's correlation test (two-sided). Source data are provided as a Source Data file.

Fig. 2 | Evidence of goal and uncertainty representation

All data are presented as mean \pm SEM from $n = 20$ participants. Asterisks denote statistical significance (paired t-test against the chance level, *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$). All statistical tests were two-sided. See Supplementary Table 2 for full statistical information. (a) Decoding accuracy of specific goals as a function of trial events. The x-axis labels "fix" for fixation, and "S1-S3" and "A1-A2" correspond to the states and actions as illustrated in Fig. 1a. The apostrophe (') denotes events in the subsequent trial. The event-specific neural measures are derived from fMRI data scanned in the corresponding time bin. The chance level is $\frac{1}{3}$, indicated by the dashed line. (b) Average goal decoding accuracy across the trial events (S1-fix'). (c) Decoding accuracy of uncertainty as a function of trial events. The chance level is 0.5, as indicated by the dashed line. Red asterisks denote statistical significance of uncertainty decoding in the specific goal condition, while blue asterisks indicate significance in the non-specific goal condition. (d) Average uncertainty decoding accuracy across the trial events (S1-fix'). Source data are provided as a Source Data file.

Fig. 3 | Shattering analysis for neural goal and uncertainty embeddings

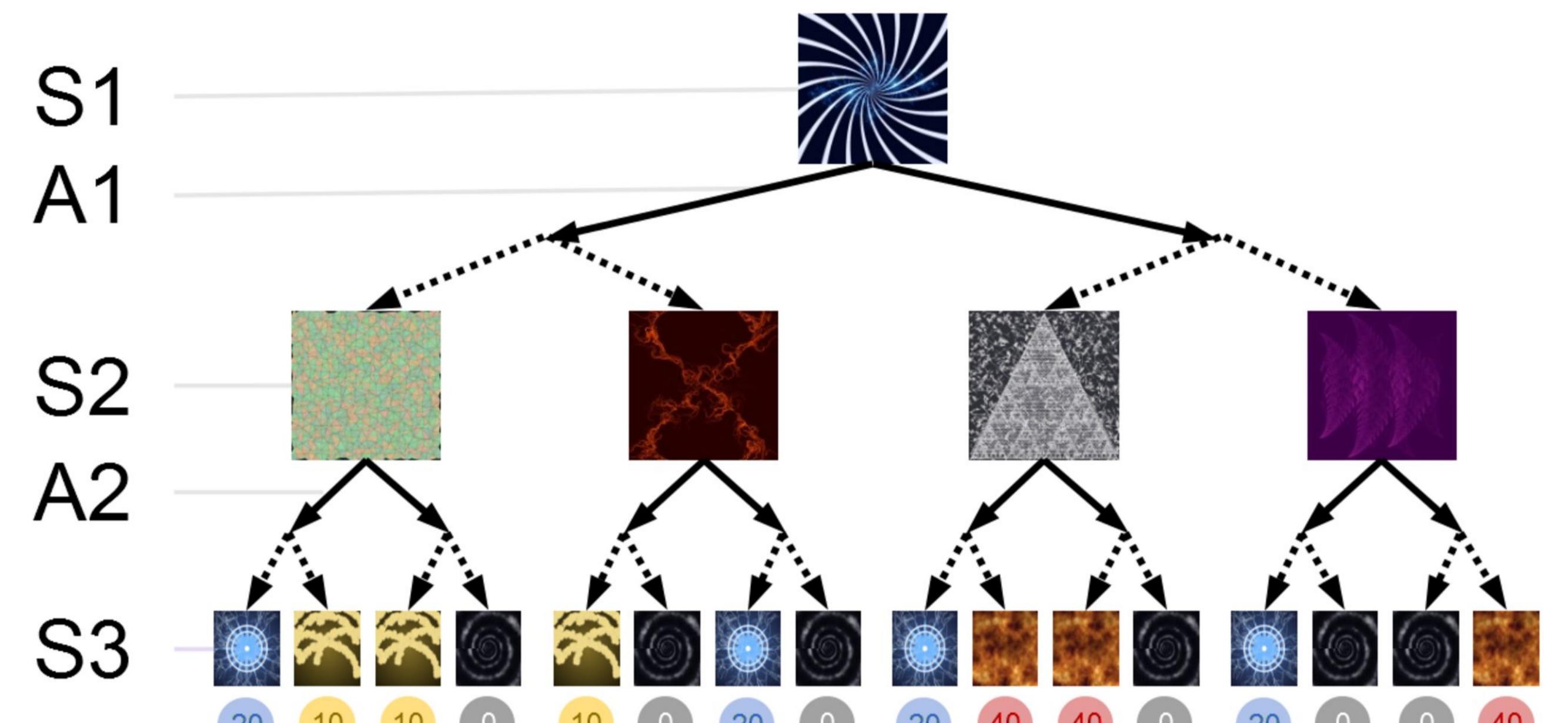
(a) Hypothetical neural embeddings for goal and uncertainty and corresponding linear separabilities. While there are three specific goals (red, blue, yellow), only two are depicted for simplicity. The full version of the class labeling and classification types for all dichotomies are in Supplementary Fig. 6. (b) SD for the different types of dichotomies. We trained separate SVM classifiers for each event (S1-fix') and averaged their decoding accuracies to obtain a single SD value. Since SD represents average binary classification accuracy, the chance level is 0.5. Statistically significant SDs were represented as color bars (Supplementary Fig. 7). Data are presented as mean \pm SEM from $n = 20$ participants. Asterisks denote statistical significance (paired t-test, *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$, ****: $p < 0.0001$). All statistical tests were two-sided. (c) Correlations between the neural goal SD and the behavioral measures. Each point represents an individual participant. Solid lines present linear regression slope and dotted lines show 95% confidence bounds of the fitted line where there are statistically significant correlations. (d) Correlation coefficients between the neural SD and the behavioral measures. Only statistically significant correlations are represented with filled bars (Pearson's correlation, *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$; two-sided test). For all the exploratory correlation analyses, multiple comparison corrections were applied with a false-discovery rate (Benjamini–Hochberg procedure) for the number of ROIs with $q = 0.05$. See Supplementary Table 3 for full statistical information. Source data are provided as a Source Data file.

Fig. 4 | Neural goal robustness and correlations with behavior

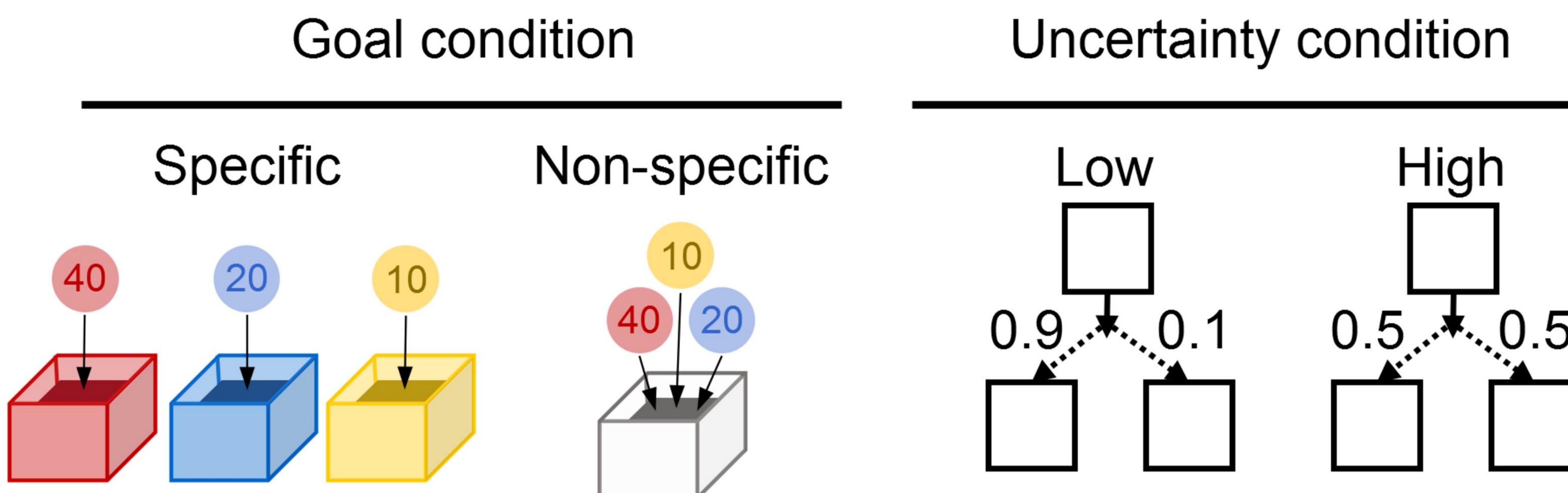
(a) Goal CCGP across different uncertainty conditions and goal SD in different uncertainty conditions. The SD values were calculated similarly to Fig. 3b, with the only difference being the division into low and high uncertainty conditions. Each CCGP or SD value represents the average value across trial events (S1-fix'). Data are presented as mean \pm SEM from $n = 20$ participants. Pairwise comparisons of CCGP, SD in low uncertainty, and SD in high uncertainty across all regions (paired t-test) showed no significant differences. All statistical tests were two-sided. (b) Correlation coefficients between the goal CCGP and the behavioral measures. Only statistically significant correlations are represented with filled bars (Pearson's correlation, *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$; two-sided test), corrected for the number of ROIs (Benjamini-Hochberg procedure, $q=0.05$). '+' indicates weak statistical significance (significant before multiple correction, +: $p < 0.05$). See Supplementary Table 4 for full statistical information. For the complete scatter plots, see Supplementary Fig. 8. Source data are provided as a Source Data file.

a

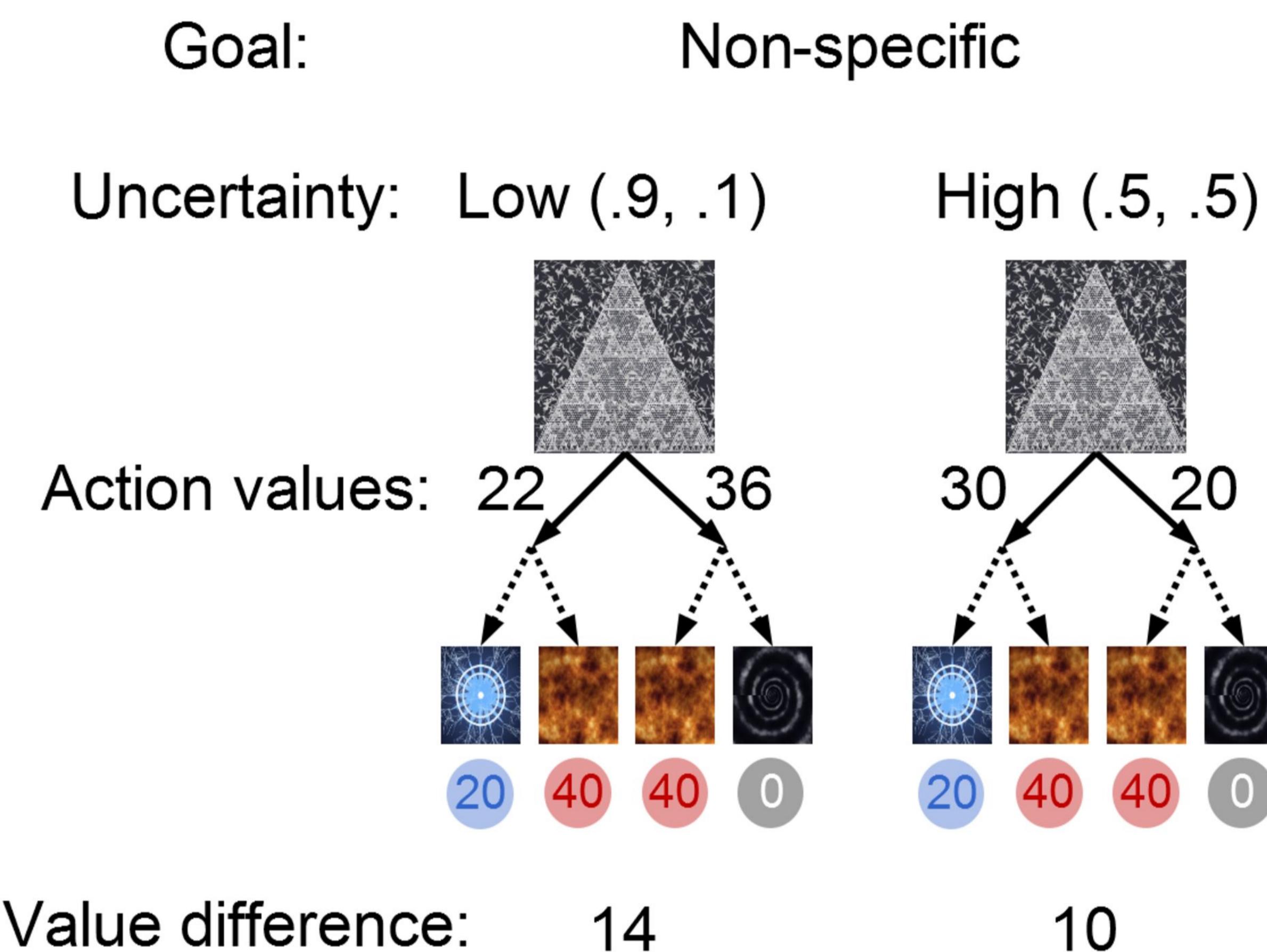
Two-stage Markov decision task

**b**

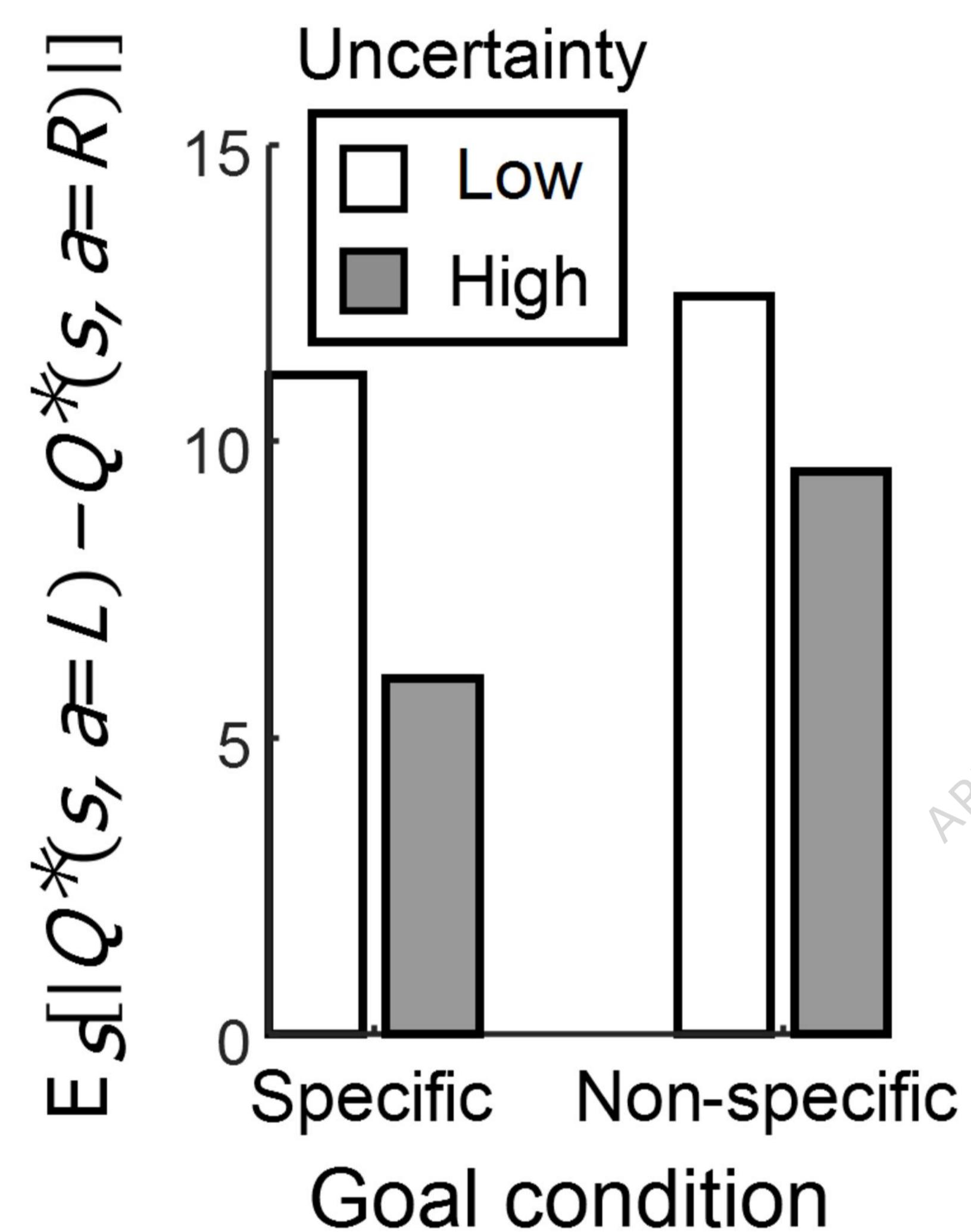
Block conditions (contexts)



c Action value computation example

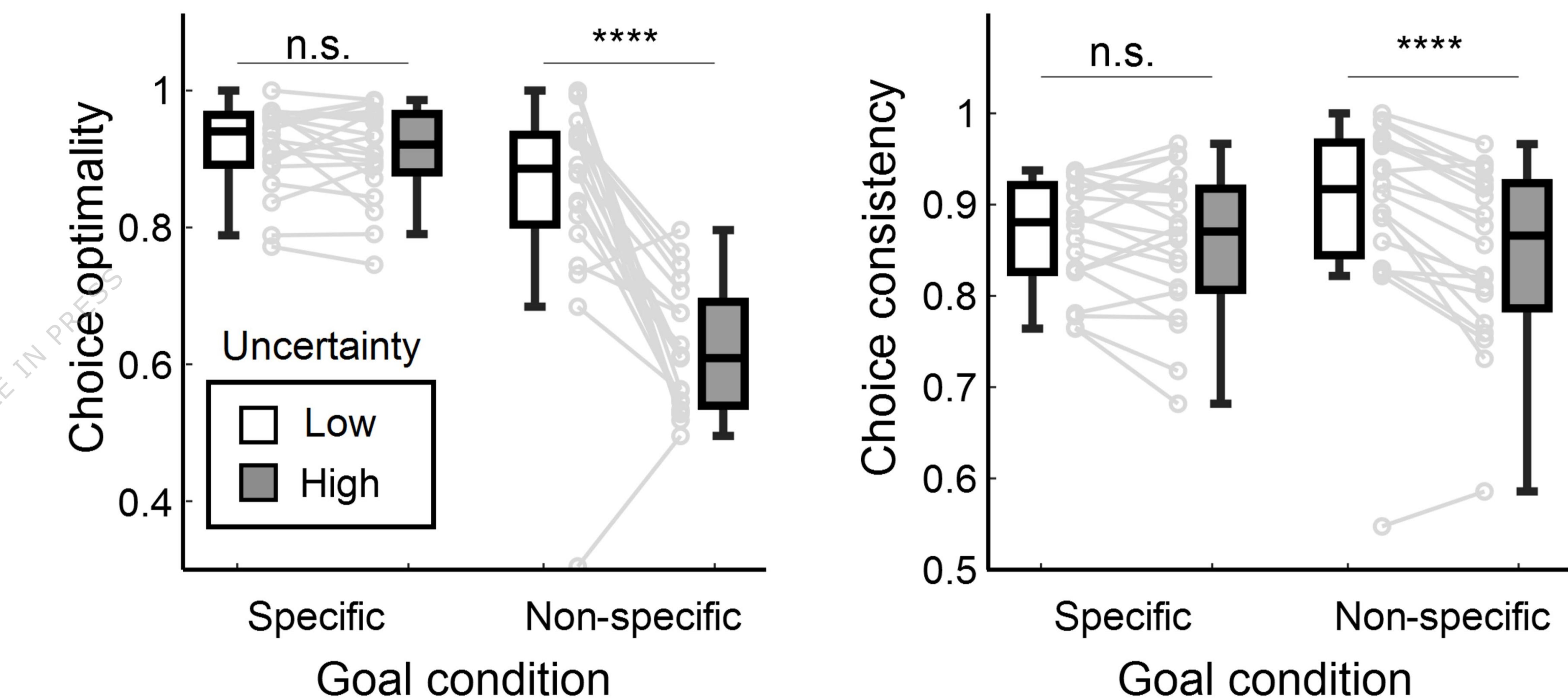


d Action value difference

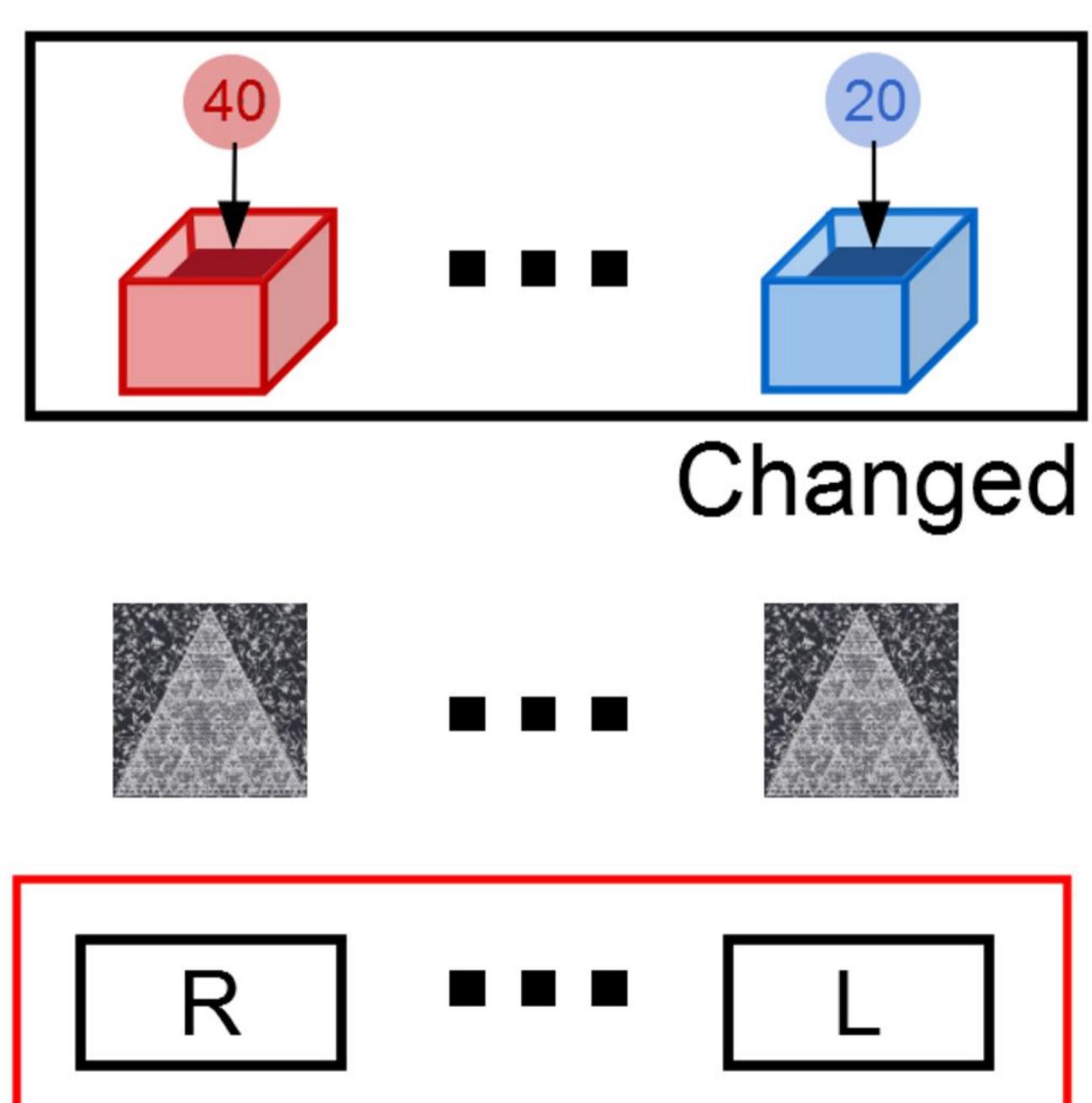


d

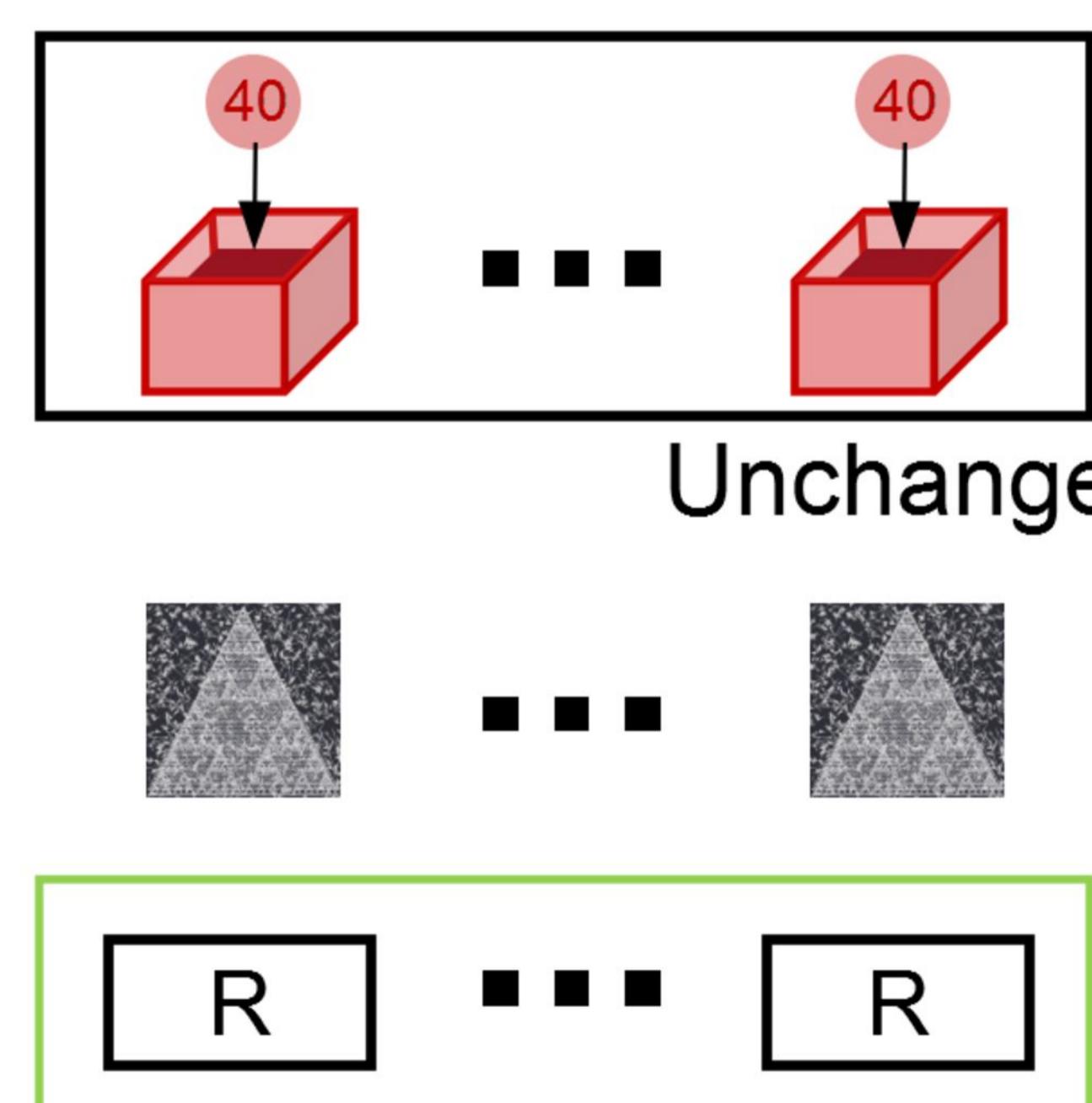
Goal x uncertainty effect on behavior

**e**

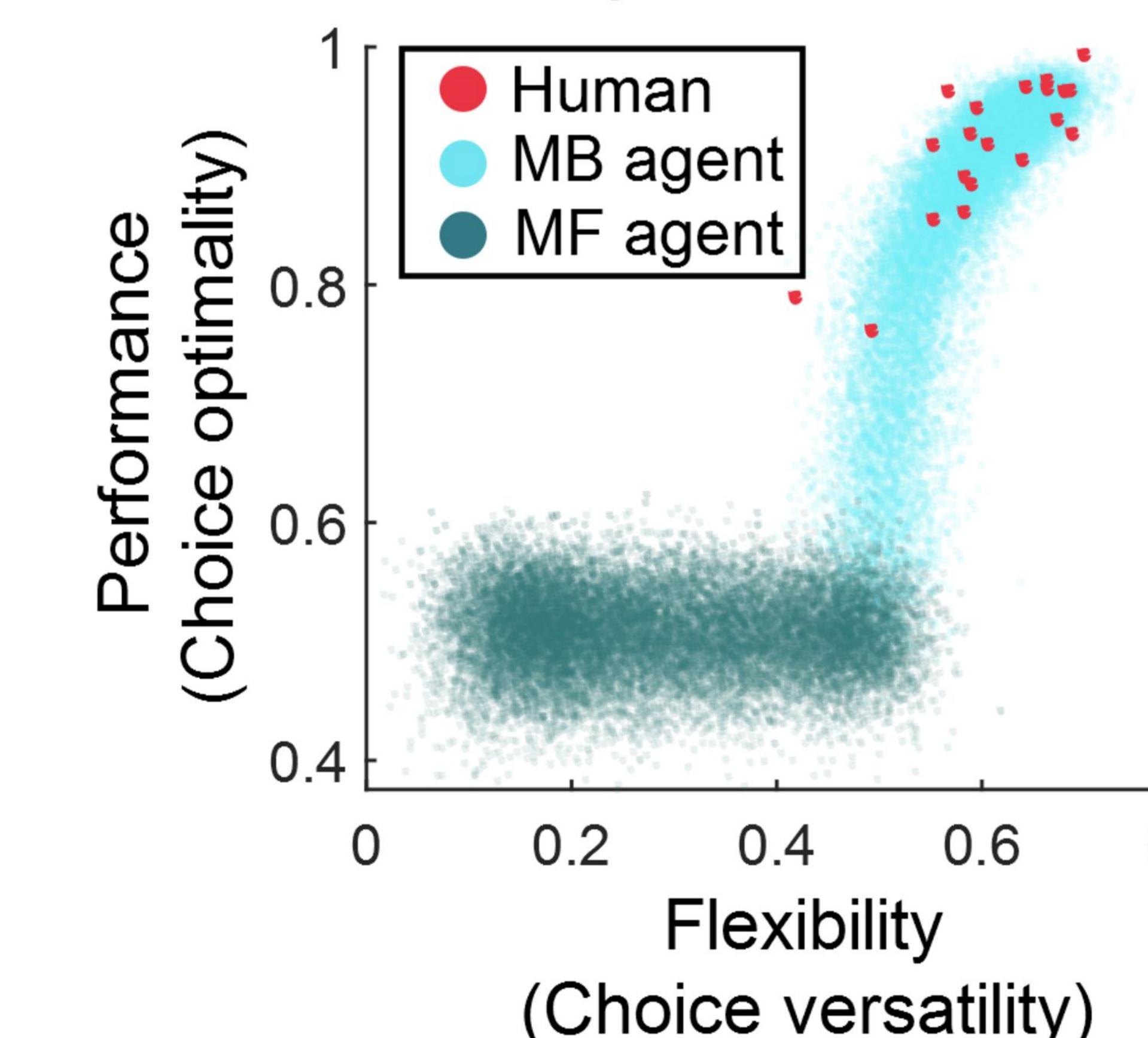
Flexibility: choice versatility



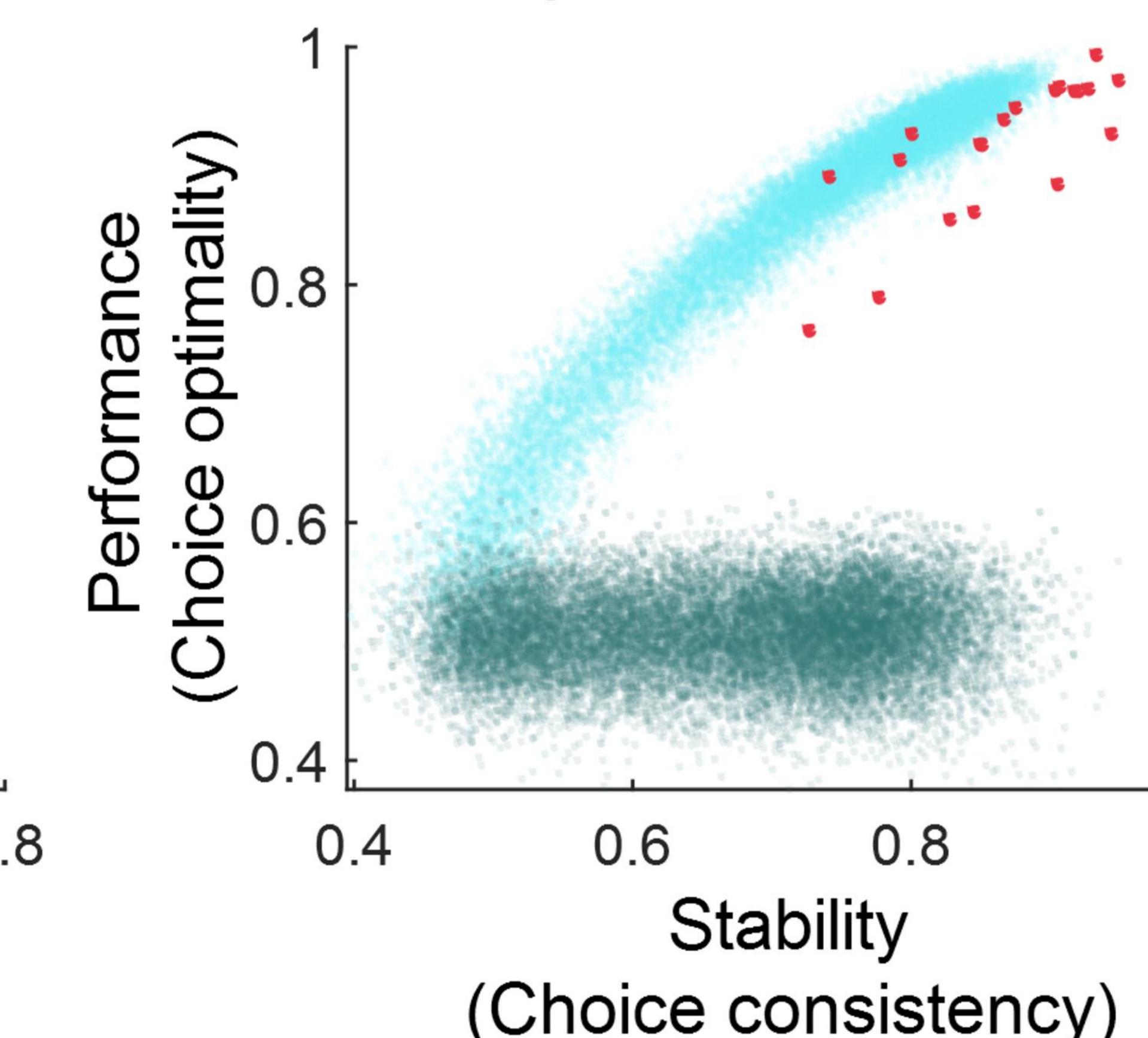
Stability: choice consistency

**f**

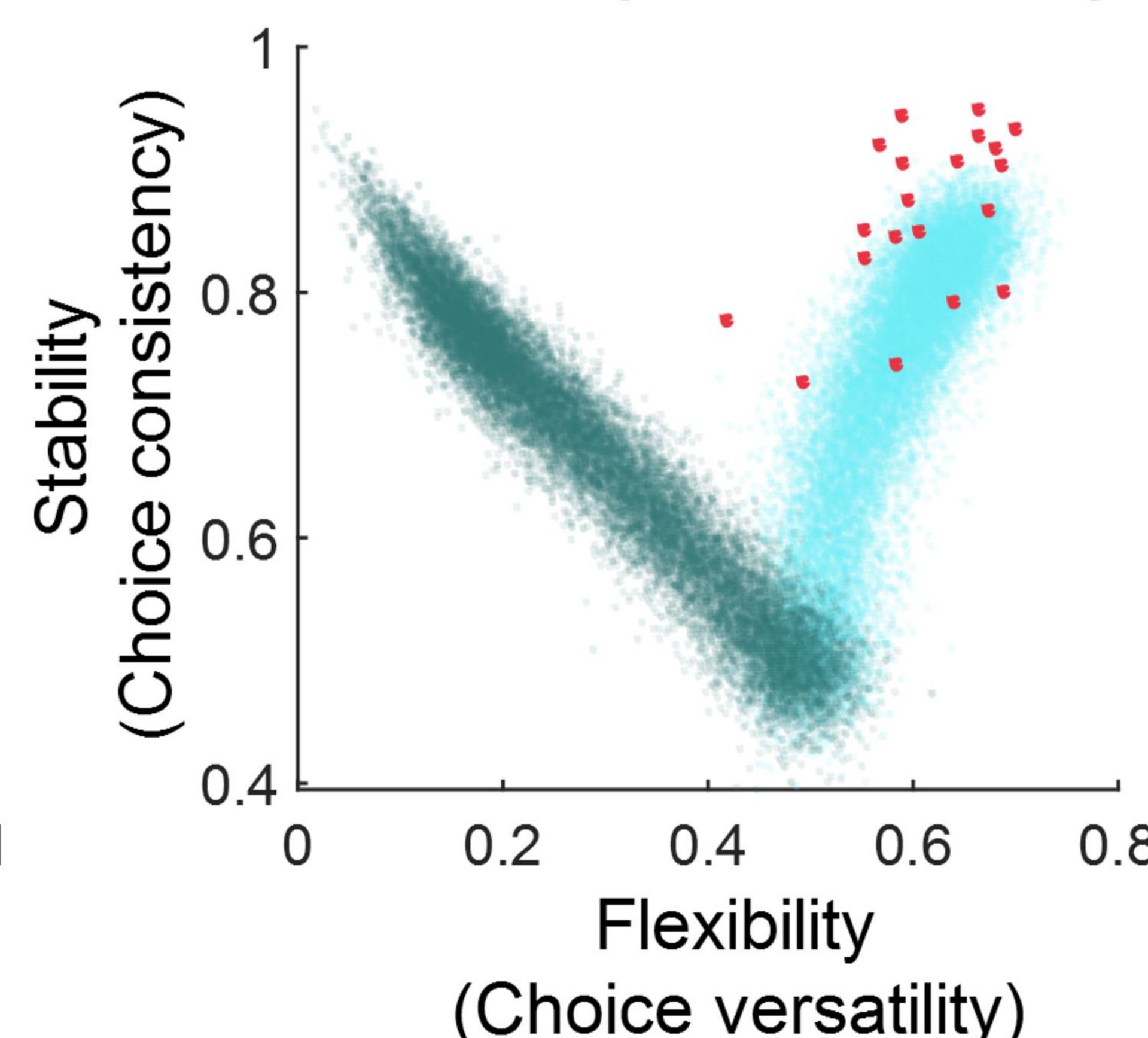
Flexibility vs Performance

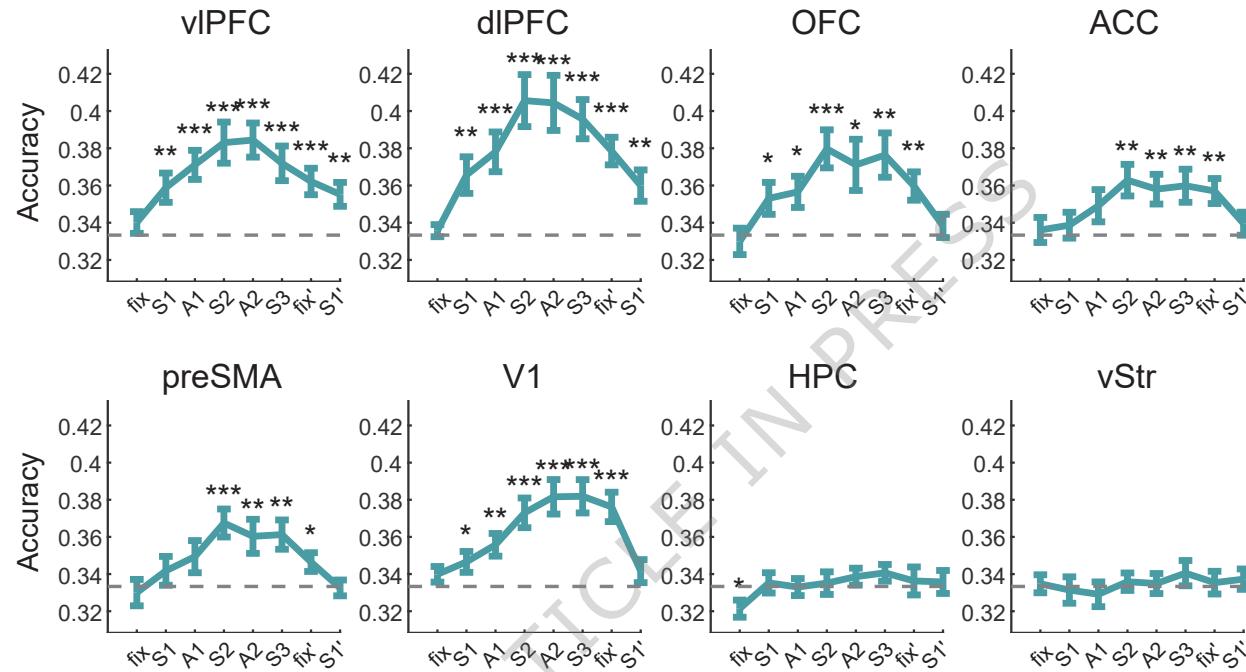
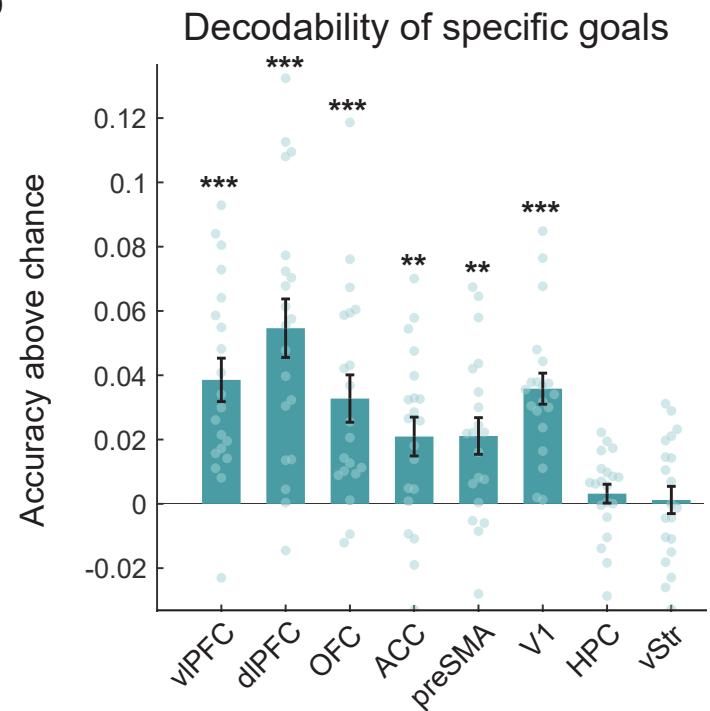
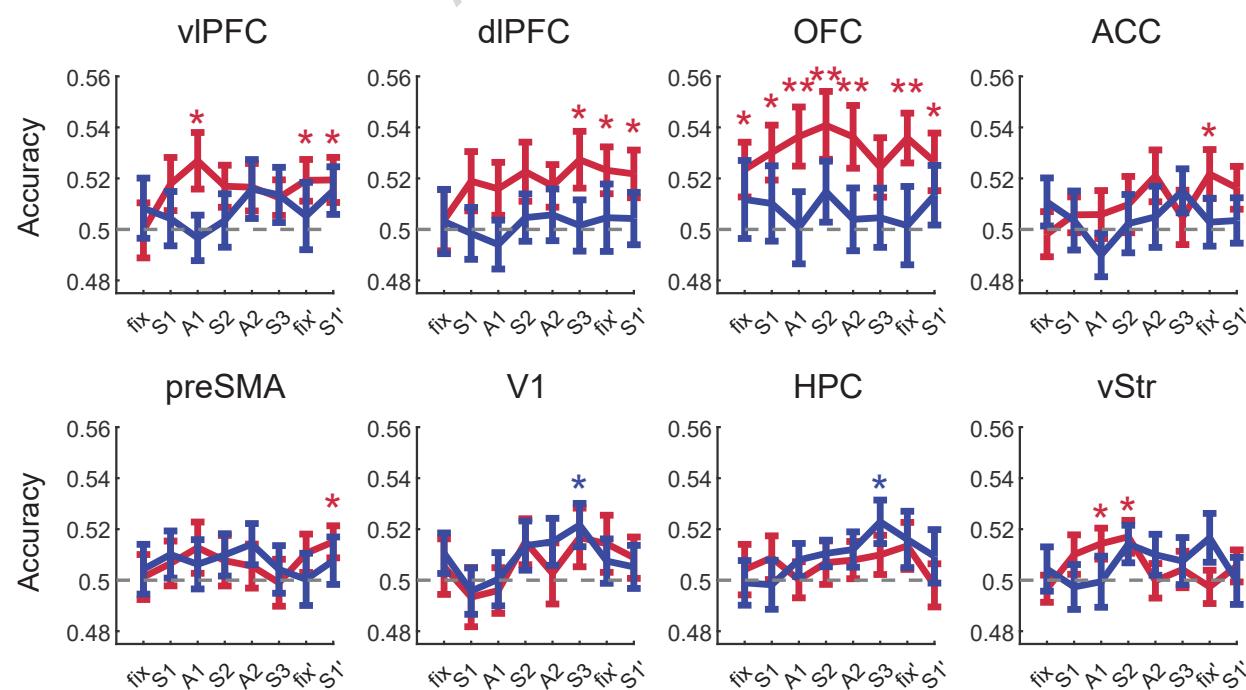
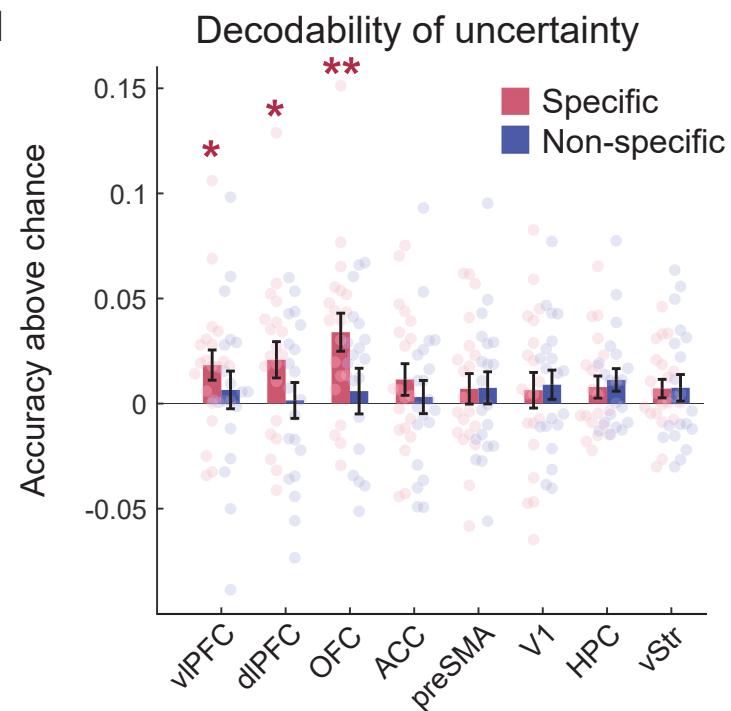


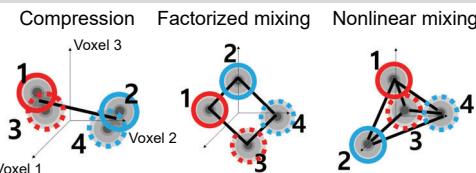
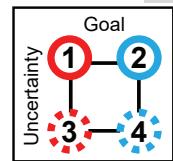
Stability vs Performance



Flexibility vs Stability



a**b****c****d**

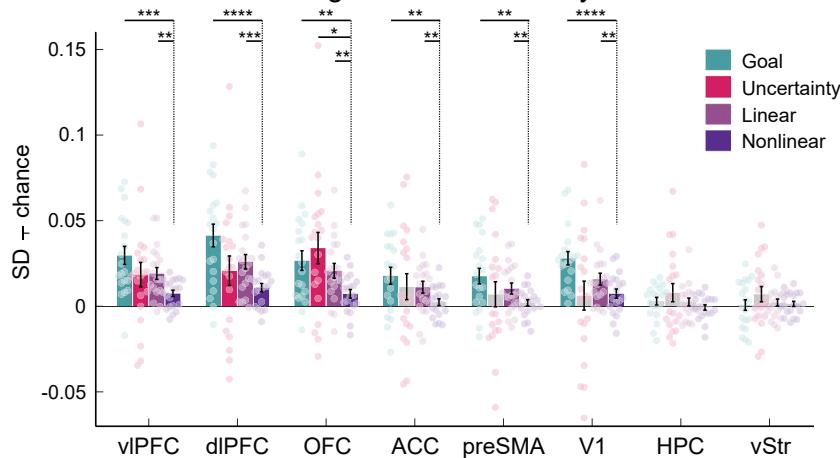
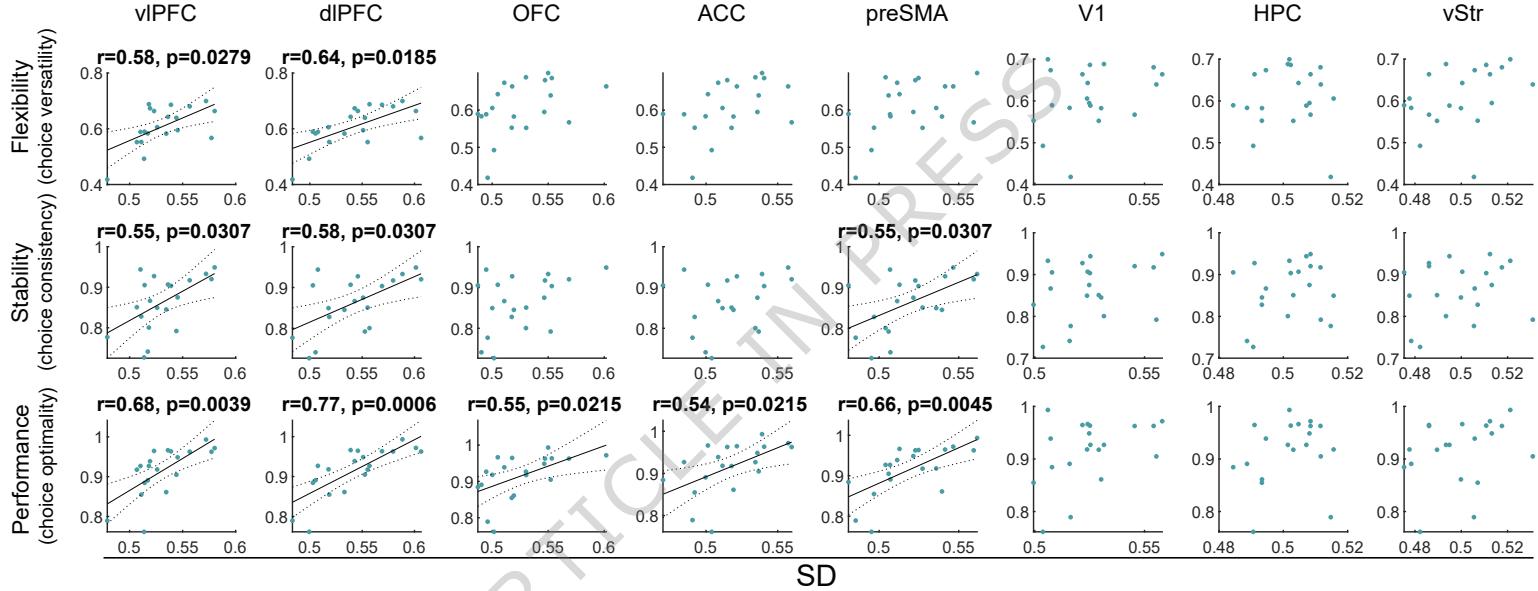
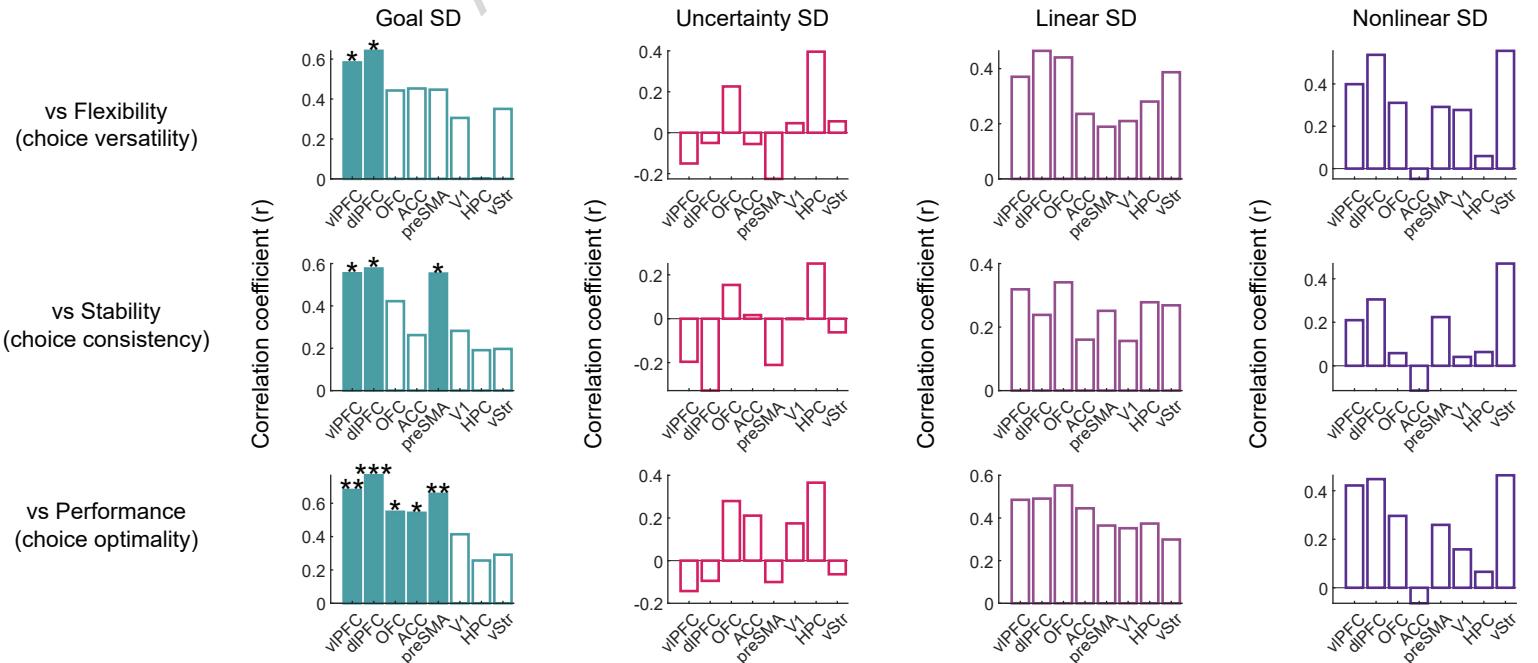
a

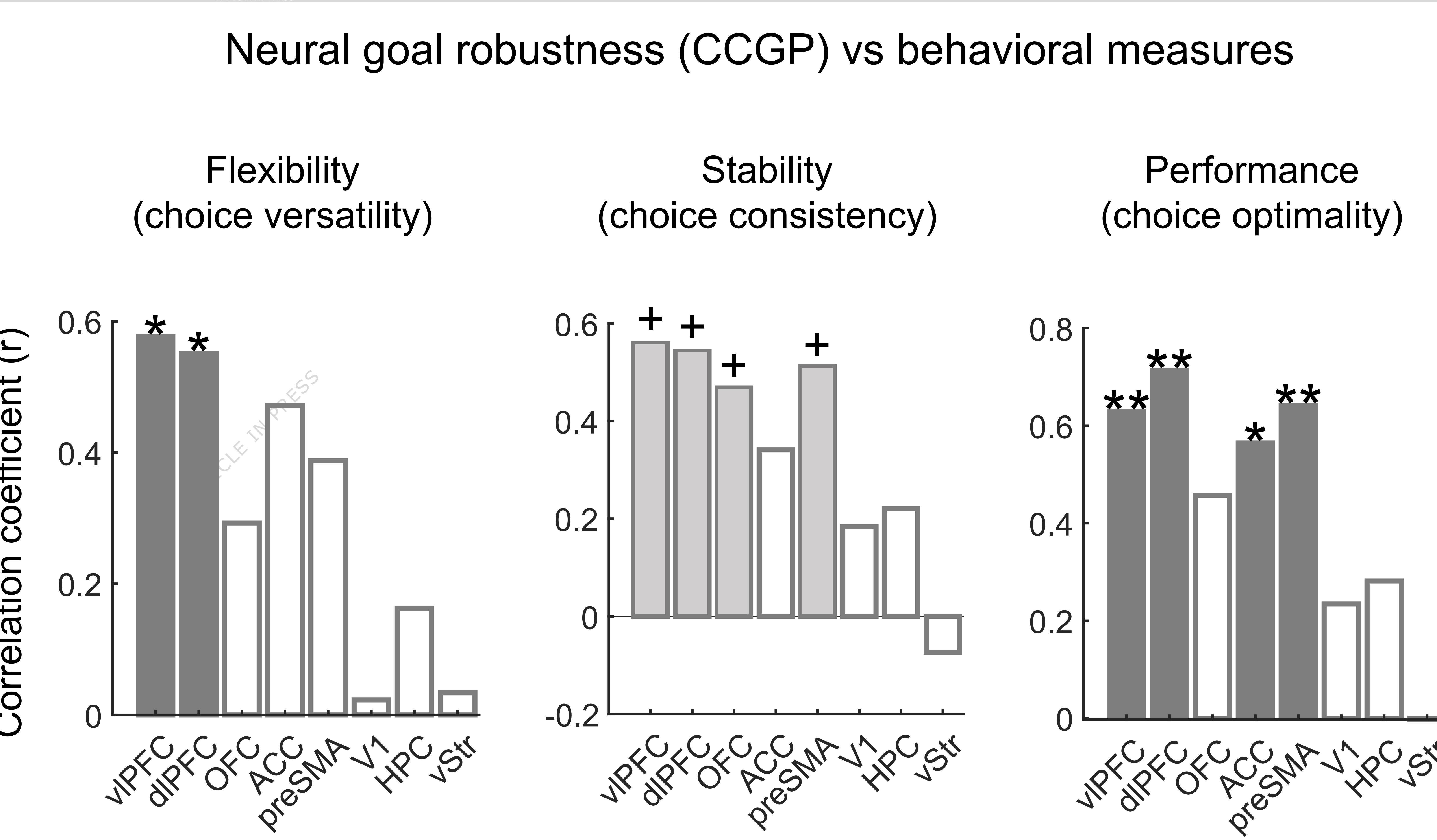
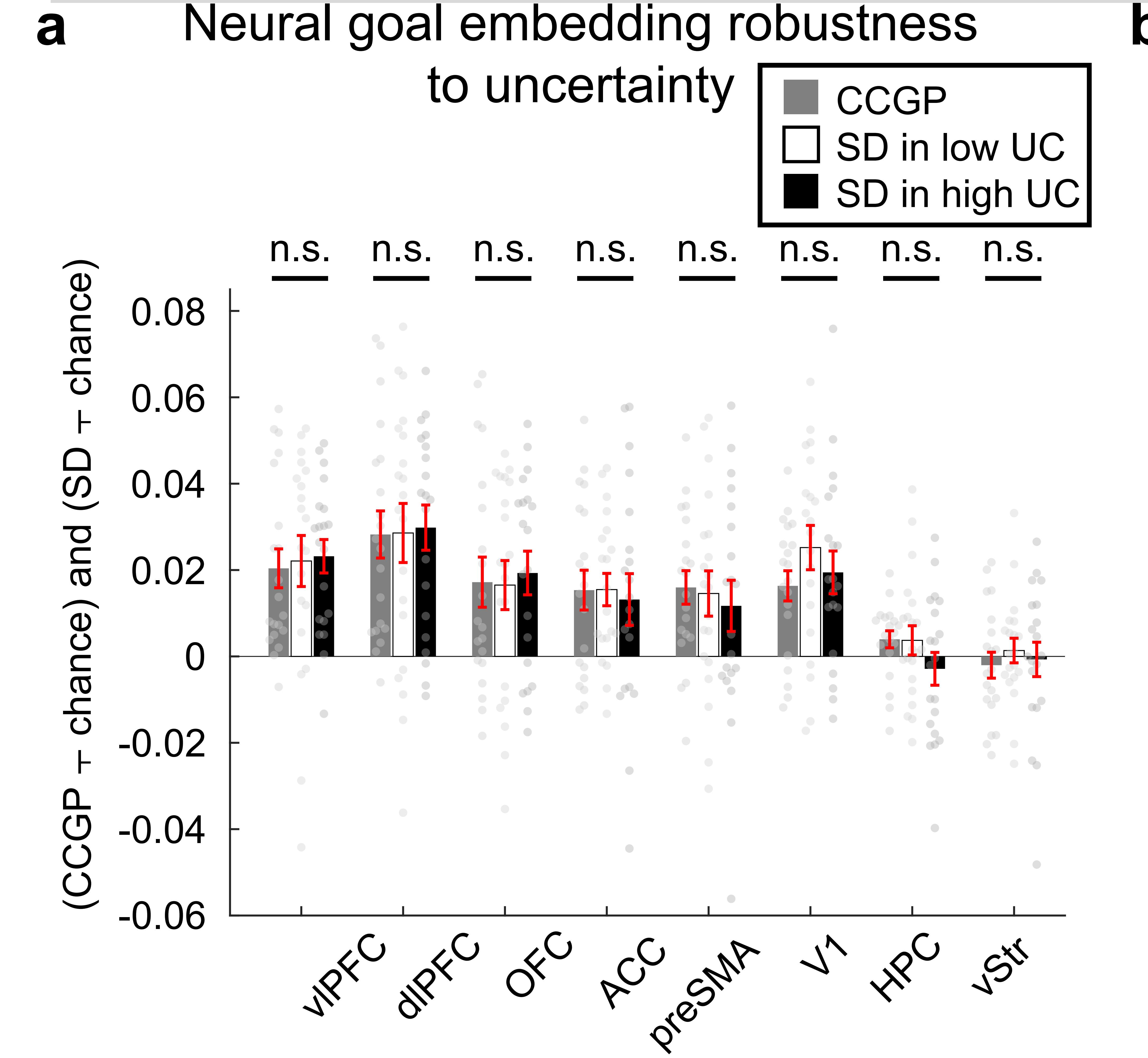
Classification type

Linear separability	
1	Goal
✓	✓
✗	✗
✗	✓
✗	✗

3 Linear

4 Nonlinear

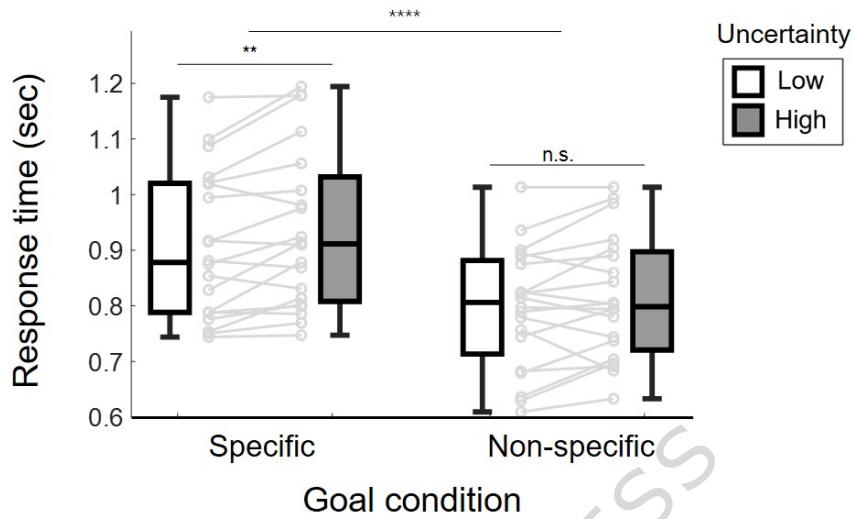
goal and uncertainty**c****Neural goal separability (SD) vs behavioral measures****d****Neural separabilities (SD) vs behavioral measures****f**



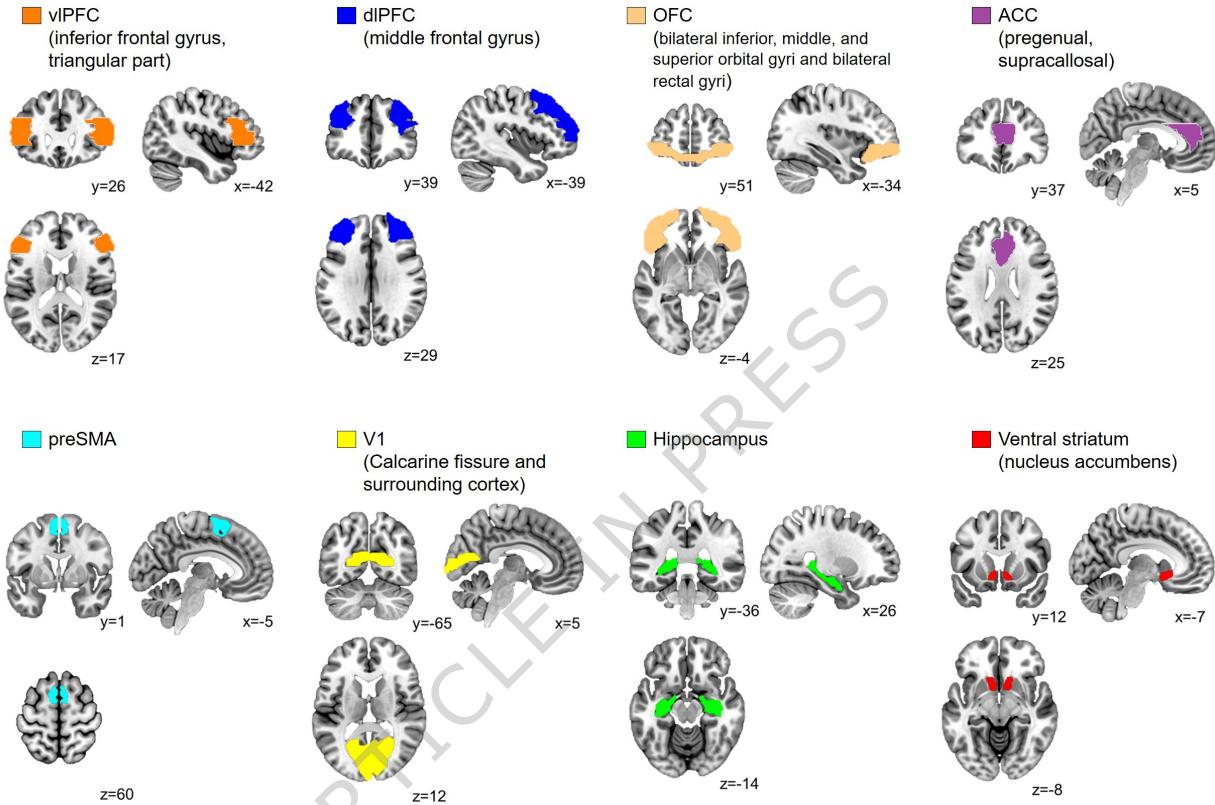
Supplementary information

Factorized embedding of goal and uncertainty in the lateral prefrontal cortex guides stably flexible learning

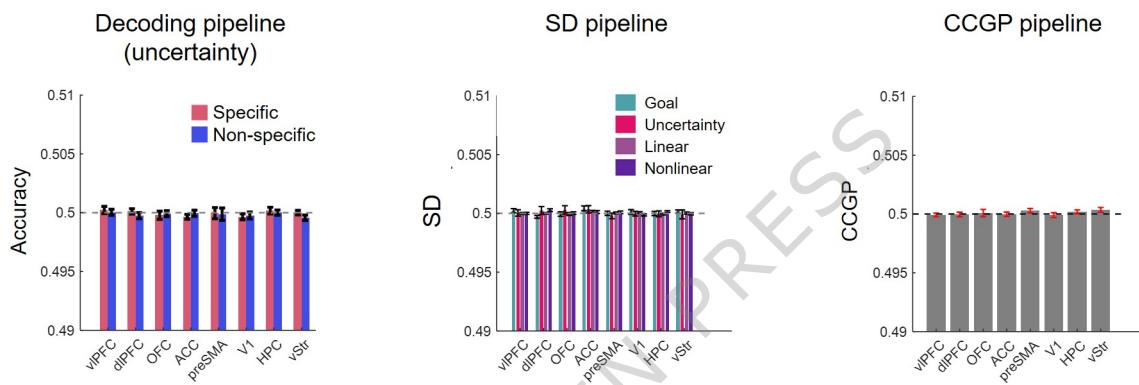
Yoondo Sung, Mattia Rigotti, and Sang Wan Lee



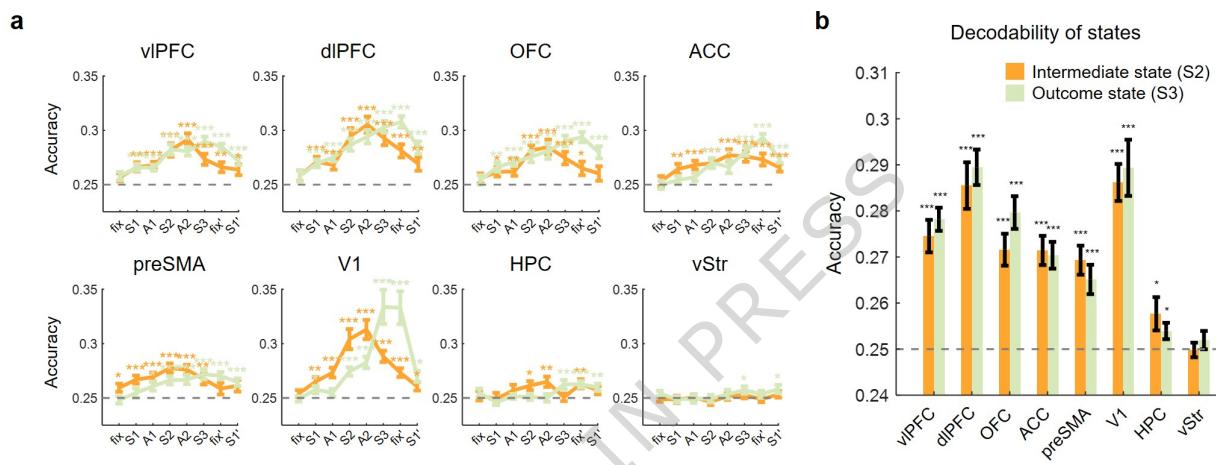
Supplementary Figure 1 | Mean response times during choice. Each dot represents one participant ($n=20$). For every experimental condition, response times from Stage 1 and Stage 2 were averaged across all trials. Mean response times were significantly longer in the specific-goal condition than in the non-specific goal condition (two-sided paired t-test, $t(19) = 5.431, p = 3.065e - 05$). Within the specific-goal condition, trials performed under high uncertainty elicited significantly longer response times than those under low uncertainty (two-sided paired t-test, $t(19) = -3.415, p = 0.0029$), whereas uncertainty had no effect in the non-specific goal condition. Source data are provided as a Source Data file.



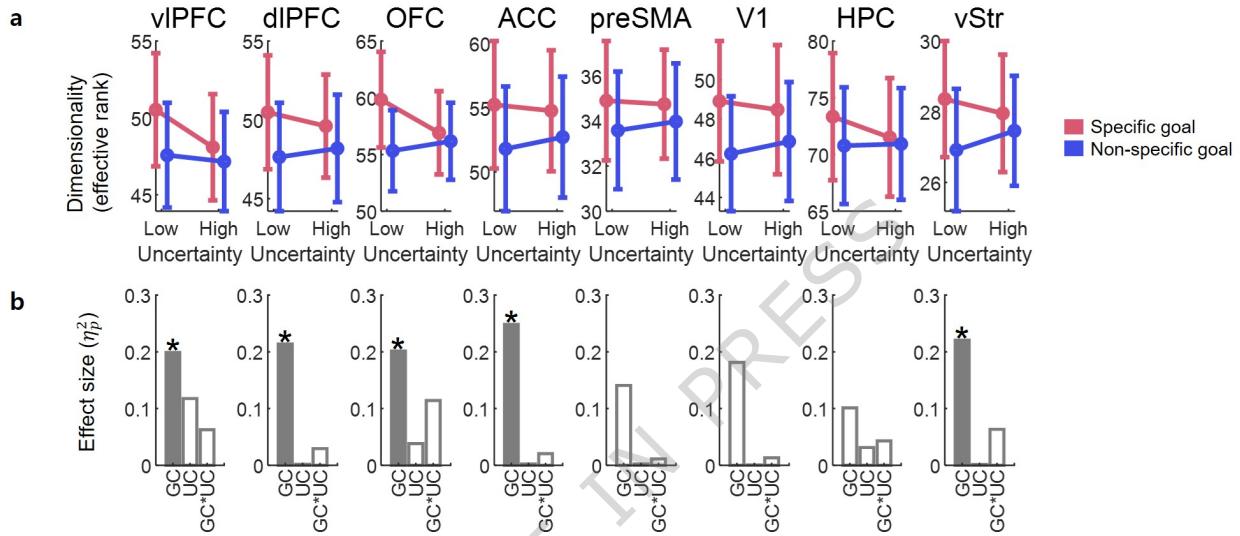
Supplementary Figure 2 | Anatomically defined ROIs. Binary masks for the eight ROIs—ventrolateral prefrontal cortex (vIPFC), dorsolateral PFC (dlPFC), orbitofrontal cortex (OFC), anterior cingulate cortex (ACC), pre-supplementary motor area (preSMA), primary visual cortex (V1), hippocampus (HPC), and ventral striatum (vStr)—are overlaid bilaterally on the MNI152 T1-weighted template (coronal, sagittal, and axial views). All masks were taken from the AAL3 atlas⁶² except the preSMA, which was obtained from the JuBrain Anatomy Toolbox⁹⁷.



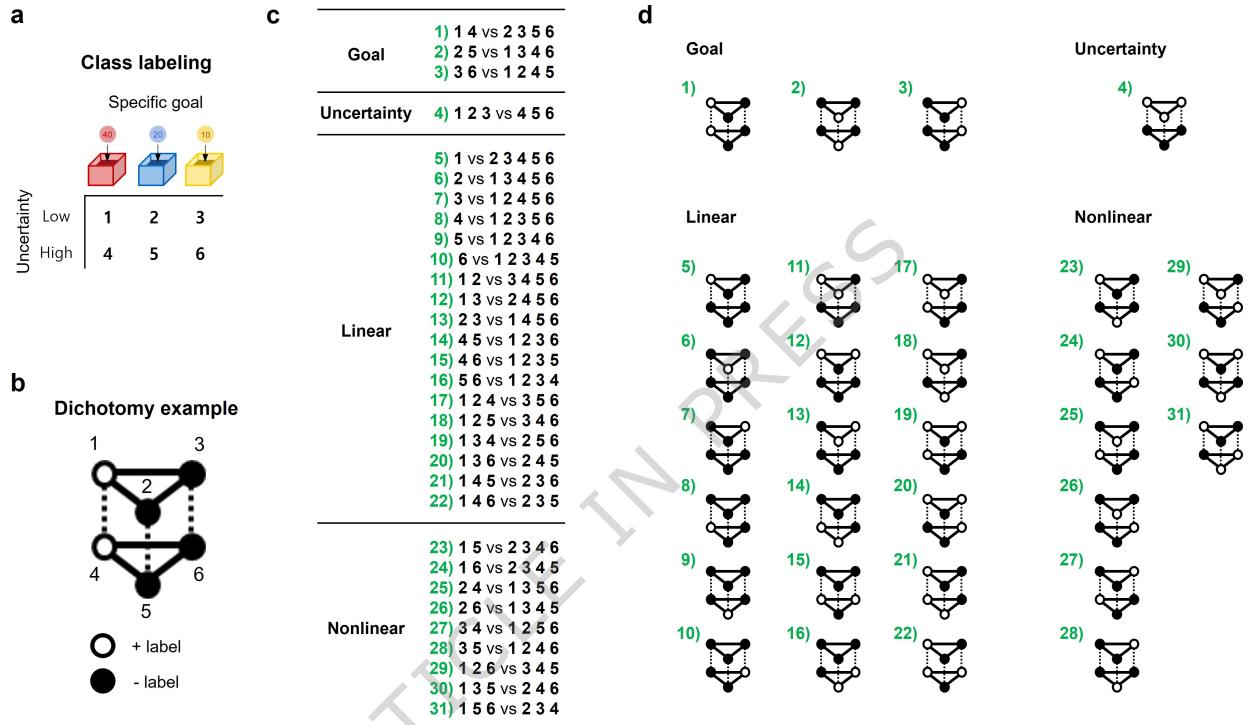
Supplementary Figure 3 | Shuffled-label controls for decoding analyses. We assessed the specificity of the three decoding analysis—standard decoding accuracy, shattering dimensionality (SD), and cross-condition generalization performance (CCGP)—by repeating each analysis with randomly permuted class labels. To correct for class imbalance, we performed 100 independent label-undersampling iterations with distinct random seeds, reporting the mean across iterations (mirroring the procedures used in Fig. 2b,d, Fig. 3b, and Fig. 4a). Under these shuffled-label conditions, every metric in every ROI remained statistically indistinguishable from chance (two-sided paired t-test vs. chance level), confirming that the main-text results cannot be attributed to spurious label structure. Error bars represent the standard error across participants. Source data are provided as a Source Data file.



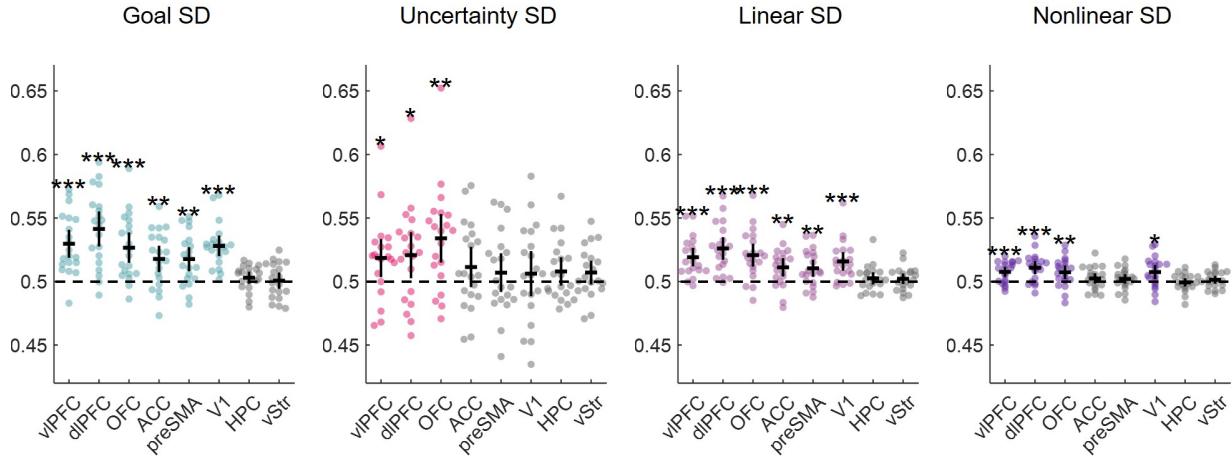
Supplementary Figure 4 | Evidence of state representation. Error bars represent the standard error across participants. Orange asterisks denote the statistical significance of the intermediate state decoding, while green asterisks indicate the significance of the outcome state (paired t-test against the chance level, *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$). All statistical tests were two-sided. (a) Decoding accuracy of state as a function of trial events. The x-axis labels "fix" for fixation, and "S1-S3" and "A1-A2" correspond to the states and actions as illustrated in Fig. 1a. The apostrophe (') denotes events in the subsequent trial. The event-specific neural measures are derived from fMRI data scanned in the corresponding time bin. The chance level is 0.25, indicated by the dashed line. (b) Average state decoding accuracy across the trial events (S1-fix'). Source data are provided as a Source Data file.



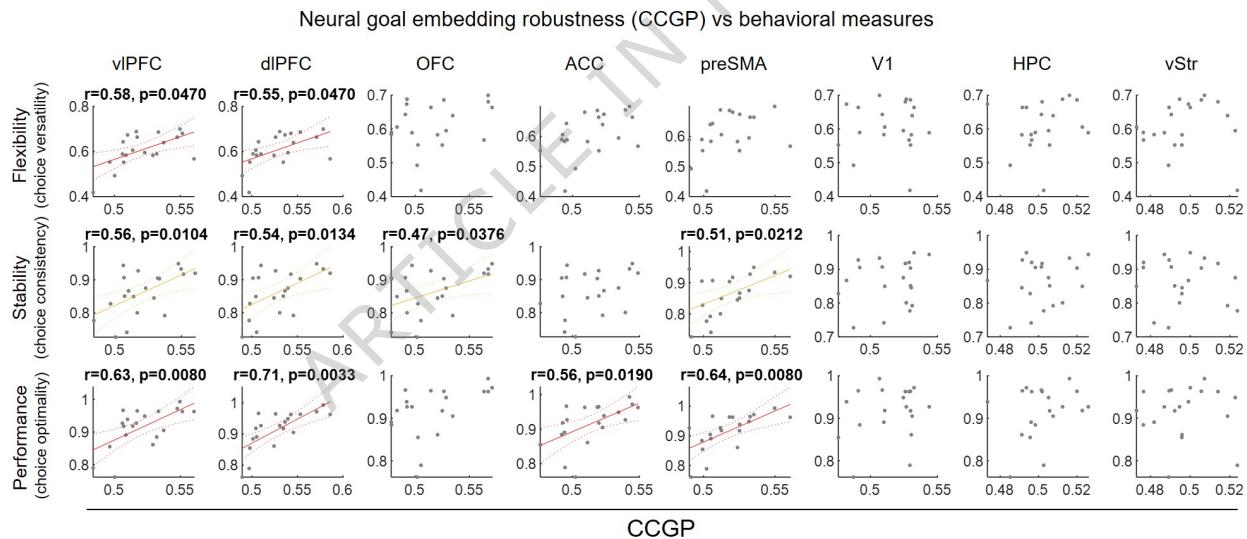
Supplementary Figure 5 | Effect of goal and uncertainty condition on the neural dimensionality of brain regions. (a) The neural dimensionality of brain regions measured under different goal and uncertainty conditions. The dimensionality was quantified using the effective rank measure¹⁰¹. This measure involves performing PCA on the BOLD signal from each brain region to obtain the eigenspectrum, calculating its entropy H , and then exponentiating the result. As defined by the previous work¹⁰¹, the effective rank $ER = \exp(H)$. Error bars represent the standard error across participants. Effective rank values were averaged across the left and right hemispheres. (b) The effect sizes obtained from a two-way repeated measures ANOVA. GC represents the goal condition, UC represents the uncertainty condition, and GC*UC denotes the interaction between these two conditions. Only significant effects are shown with filled bars (*: $p < 0.05$).



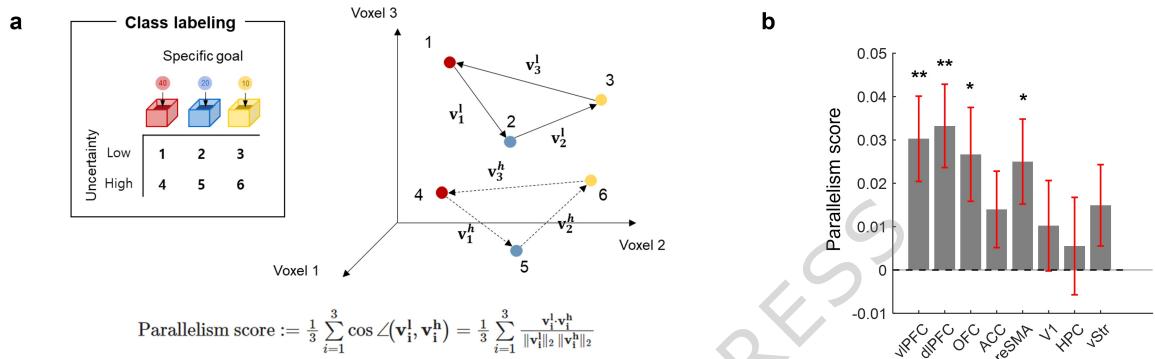
Supplementary Figure 6 | Categorization of all dichotomies. (a) The six classes determined by the specific goal and uncertainty conditions. (b) A schematic of a dichotomy based on the six classes. (c) The four categories (goal, uncertainty, linear, and nonlinear) for all dichotomies. There are 2^6 ways of binary labeling with the six classes. The actual number of dichotomies reduces to $\frac{2^6-2}{2}$ by excepting the two cases of all positive or negative labeling and removing half of the duplicated cases due to the symmetry of binary labeling. (d) Visualization of all categorized dichotomies.



Supplementary Figure 7 | Statistical significance of SD for the different types of dichotomies The chance levels are indicated by the dashed lines. Statistically significant SD results are shown in the respective type colors, while non-significant results are shown in gray (paired t-test against the chance level, *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$). Bold bars represent the mean and the 95% confidence interval.



Supplementary Figure 8 | Correlations between goal CCGP and behavioral measures. Each point represents an individual participant. Solid lines represent the linear regression slope, and dotted lines show the 95% confidence bounds of the fitted line where there are statistically significant correlations. The p-values associated with the red regression slopes have been adjusted for multiple comparisons across the eight ROIs using the Benjamini–Hochberg false-discovery-rate procedure ($q = 0.05$). The yellow values show the corresponding nominal (uncorrected) p-values; none of these remained significant after correction.



Supplementary Figure 9 | Parallelism score results. (a) To determine whether goal representations are arranged similarly across uncertainty levels, we computed the parallelism score (PS) at the level of mean multivoxel patterns. Mean patterns of the six classes (left) are displayed as coloured points in an illustrative multivoxel activity space (right). Solid arrows trace the coding directions that connect specific-goal pairs under low uncertainty, whereas dashed arrows depict the corresponding directions under high uncertainty. For each goal pair, PS is the cosine similarity between the two coding vectors ($\mathbf{v}_i^l, \mathbf{v}_i^h$) and the final score is their mean. To remain consistent with the other validated fMRI analyses in the main text, PS values were averaged across scanning sessions. Class imbalance was controlled by performing 100 random undersampling iterations whenever class means were estimated. (b) PSs were significantly greater than zero only in vlPFC, dlPFC, OFC, and preSMA (two-sided one-sample t-test; *: $p < 0.05$ and **: $p < 0.01$). Error bars indicate the between-subject standard error of the mean ($n = 20$). Consistent with the shattering- and CCGP-analysis procedures, we computed the PS separately for each trial event (S1-fix') and report the mean across events. Source data are provided as a Source Data file.

Supplementary Table 1 | Statistical details for Fig. 1

Figure panel	Condition	Description	n	Test	Statistic	p-value
1d, left	Specific goal condition	Trial-averaged choice optimality comparison between low vs high uncertainty conditions across participants	20	Paired t-test (two-sided)	t(19)=0.887	p=3.862e-01 (n.s.)
	Non-specific goal condition		20		t(19)=5.975	p=9.463e-06 (****)
1d, right	Specific goal condition	Trial-averaged choice consistency comparison between low vs high uncertainty conditions across participants	20	Paired t-test (two-sided)	t(19)=0.891	p=3.842e-01 (n.s.)
	Non-specific goal condition		20		t(19)=5.294	p=4.141e-05 (****)
1f, left	Human	Correlation between trial-averaged choice versatility and choice optimality across human participants	20	Pearson's correlation	r=0.827, t(18)=6.244	p=6.859e-06 (****)
	Model-based agent	Correlation between trial-averaged choice versatility and choice optimality across virtual agents (distinct random seeds)	20000		r=0.799, t(19998)=188.200	p=0.000e+00 (****)
	Model-free agent		20000		r=-0.035, t(19998)=-4.886	p=1.037e-06 (****)
1f, middle	Human	Correlation between trial-averaged choice consistency and choice optimality across human participants	20	Pearson's correlation	r=0.779, t(18)=5.269	p=5.208e-05 (****)
	Model-based agent	Correlation between trial-averaged choice consistency and choice optimality across virtual agents (distinct random seeds)	20000		r=0.963, t(19998)=506.147	p=0.000e+00 (****)
	Model-free agent		20000		r=0.042, t(19998)=5.989	p=2.154e-09 (****)
1f, right	Human	Correlation between trial-averaged choice versatility and choice consistency across participants	20	Pearson's correlation	r=0.541, t(18)=2.726	p=1.386e-02 (*)
	Model-based agent	Correlation between trial-averaged choice versatility and choice consistency across agents (distinct random seeds)	20000		r=0.829, t(19998)=209.498	p=0.000e+00 (****)
	Model-free agent		20000		r=-0.964, t(19998)=-516.447	p=0.000e+00 (****)

Supplementary Table 2 | Statistical details for Fig. 2

Figure panel	Condition	Description	n	Test	Statistic	p-value
2a	vIPFC	Comparison between event-specific decoding accuracy against the chance level (1/3)	20	Paired t-test (two-sided)	fix : $t(19)=1.164$	$p=2.587e-01$ (n.s.)
					S1 : $t(19)=3.234$	$p=4.371e-03$ (**)
					A1 : $t(19)=4.825$	$p=1.175e-04$ (***)
					S2 : $t(19)=4.494$	$p=2.486e-04$ (***)
					A2 : $t(19)=5.571$	$p=2.256e-05$ (***)
					S3 : $t(19)=4.145$	$p=5.506e-04$ (***)
					fix' : $t(19)=4.050$	$p=6.842e-04$ (***)
					S1' : $t(19)=3.413$	$p=2.918e-03$ (**)
2a	dIPFC	Comparison between event-specific decoding accuracy against the chance level (1/3)	20	Paired t-test (two-sided)	fix : $t(19)=0.736$	$p=4.707e-01$ (n.s.)
					S1 : $t(19)=3.278$	$p=3.957e-03$ (**)
					A1 : $t(19)=4.188$	$p=4.989e-04$ (***)
					S2 : $t(19)=5.175$	$p=5.385e-05$ (***)
					A2 : $t(19)=4.798$	$p=1.250e-04$ (***)
					S3 : $t(19)=5.914$	$p=1.079e-05$ (***)
					fix' : $t(19)=6.088$	$p=7.454e-06$ (***)
					S1' : $t(19)=3.155$	$p=5.209e-03$ (**)
2a	OFC	Comparison between event-specific decoding accuracy against the chance level (1/3)	20	Paired t-test (two-sided)	fix : $t(19)=-0.473$	$p=6.414e-01$ (n.s.)
					S1 : $t(19)=2.302$	$p=3.284e-02$ (*)
					A1 : $t(19)=2.773$	$p=1.212e-02$ (*)
					S2 : $t(19)=4.532$	$p=2.279e-04$ (***)
					A2 : $t(19)=2.743$	$p=1.294e-02$ (*)
					S3 : $t(19)=3.611$	$p=1.861e-03$ (**)
					fix' : $t(19)=3.462$	$p=2.610e-03$ (**)
					S1' : $t(19)=0.794$	$p=4.372e-01$ (n.s.)
2a	ACC	Comparison between event-specific decoding accuracy against the chance level (1/3)	20	Paired t-test (two-sided)	fix : $t(19)=0.410$	$p=6.865e-01$ (n.s.)
					S1 : $t(19)=0.767$	$p=4.524e-01$ (n.s.)
					A1 : $t(19)=1.825$	$p=8.378e-02$ (n.s.)
					S2 : $t(19)=3.481$	$p=2.499e-03$ (**)
					A2 : $t(19)=3.091$	$p=6.017e-03$ (**)
					S3 : $t(19)=2.983$	$p=7.643e-03$ (**)

				fix` : t(19)=3.525	p=2.261e-03 (**)
				S1` : t(19)=0.993	p=3.334e-01 (n.s.)
preSMA	Comparison between event-specific decoding accuracy against the chance level (1/3)	20	Paired t-test (two-sided)	fix : t(19)=-0.466	p=6.463e-01 (n.s.)
				S1 : t(19)=1.085	p=2.916e-01 (n.s.)
				A1 : t(19)=1.853	p=7.953e-02 (n.s.)
				S2 : t(19)=4.537	p=2.253e-04 (***)
				A2 : t(19)=2.942	p=8.376e-03 (**)
				S3 : t(19)=3.542	p=2.180e-03 (**)
				fix` : t(19)=2.596	p=1.773e-02 (*)
				S1` : t(19)=-0.184	p=8.563e-01 (n.s.)
V1	Comparison between event-specific decoding accuracy against the chance level (1/3)	20	Paired t-test (two-sided)	fix : t(19)=1.553	p=1.368e-01 (n.s.)
				S1 : t(19)=2.348	p=2.985e-02 (*)
				A1 : t(19)=3.672	p=1.618e-03 (**)
				S2 : t(19)=4.958	p=8.724e-05 (***)
				A2 : t(19)=5.196	p=5.140e-05 (***)
				S3 : t(19)=5.409	p=3.217e-05 (***)
				fix` : t(19)=5.413	p=3.187e-05 (***)
				S1` : t(19)=1.331	p=1.990e-01 (n.s.)
HPC	Comparison between event-specific decoding accuracy against the chance level (1/3)	20	Paired t-test (two-sided)	fix : t(19)=-2.587	p=1.809e-02 (*)
				S1 : t(19)=0.353	p=7.283e-01 (n.s.)
				A1 : t(19)=-0.069	p=9.459e-01 (n.s.)
				S2 : t(19)=0.294	p=7.722e-01 (n.s.)
				A2 : t(19)=1.152	p=2.638e-01 (n.s.)
				S3 : t(19)=1.631	p=1.193e-01 (n.s.)
				fix` : t(19)=0.394	p=6.980e-01 (n.s.)
				S1` : t(19)=0.390	p=7.006e-01 (n.s.)
vStr	Comparison between event-specific decoding accuracy against the chance level (1/3)	20	Paired t-test (two-sided)	fix : t(19)=0.298	p=7.688e-01 (n.s.)
				S1 : t(19)=-0.270	p=7.898e-01 (n.s.)
				A1 : t(19)=-0.632	p=5.349e-01 (n.s.)
				S2 : t(19)=0.542	p=5.941e-01 (n.s.)
				A2 : t(19)=0.294	p=7.721e-01 (n.s.)
				S3 : t(19)=1.053	p=3.054e-01 (n.s.)

					fix` : t(19)=0.335	p=7.409e-01 (n.s.)
					S1` : t(19)=0.716	p=4.827e-01 (n.s.)
2b	-	Comparison between event-averaged decoding accuracy against the chance level (1/3)	20	Paired t-test (two-sided)	vIPFC: t(19)=5.712	p=1.664e-05 (****)
					dIPFC: t(19)=5.999	p=8.991e-06 (****)
					OFC: t(19)=4.443	p=2.788e-04 (***)
					ACC: t(19)=3.459	p=2.627e-03 (**)
					preSMA: t(19)=3.701	p=1.518e-03 (**)
					V1: t(19)=7.401	p=5.220e-07 (****)
					HPC: t(19)=1.083	p=2.926e-01 (n.s.)
					vStr: t(19)=0.284	p=7.795e-01 (n.s.)
2c, Specific goal condition (red)	vIPFC	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)	fix : t(19)=-0.035	p=9.726e-01 (n.s.)
					S1 : t(19)=1.708	p=1.039e-01 (n.s.)
					A1 : t(19)=2.427	p=2.536e-02 (*)
					S2 : t(19)=2.053	p=5.408e-02 (n.s.)
					A2 : t(19)=1.778	p=9.138e-02 (n.s.)
					S3 : t(19)=1.789	p=8.963e-02 (n.s.)
					fix` : t(19)=2.351	p=2.968e-02 (*)
					S1` : t(19)=2.194	p=4.087e-02 (*)
					fix : t(19)=0.301	p=7.668e-01 (n.s.)
					S1 : t(19)=1.628	p=1.201e-01 (n.s.)
dIPFC	dIPFC	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)	A1 : t(19)=1.533	p=1.417e-01 (n.s.)
					S2 : t(19)=1.952	p=6.580e-02 (n.s.)
					A2 : t(19)=2.044	p=5.508e-02 (n.s.)
					S3 : t(19)=2.454	p=2.394e-02 (*)
					fix` : t(19)=2.522	p=2.076e-02 (*)
					S1` : t(19)=2.330	p=3.097e-02 (*)
					fix : t(19)=2.185	p=4.160e-02 (*)
					S1 : t(19)=2.796	p=1.153e-02 (*)
					A1 : t(19)=3.142	p=5.371e-03 (**)
					S2 : t(19)=3.069	p=6.311e-03 (**)
OFC	OFC	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)	A2 : t(19)=2.949	p=8.247e-03 (**)
					S3 : t(19)=2.085	p=5.081e-02 (n.s.)

				fix` : t(19)=3.666	p=1.643e-03 (**)
				S1` : t(19)=2.344	p=3.011e-02 (*)
ACC	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)	fix : t(19)=-0.213	p=8.340e-01 (n.s.)
				S1 : t(19)=0.813	p=4.264e-01 (n.s.)
				A1 : t(19)=0.621	p=5.419e-01 (n.s.)
				S2 : t(19)=0.883	p=3.885e-01 (n.s.)
				A2 : t(19)=2.068	p=5.254e-02 (n.s.)
				S3 : t(19)=0.442	p=6.636e-01 (n.s.)
				fix` : t(19)=2.258	p=3.592e-02 (*)
				S1` : t(19)=1.942	p=6.707e-02 (n.s.)
preSMA	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)	fix : t(19)=0.141	p=8.897e-01 (n.s.)
				S1 : t(19)=0.748	p=4.633e-01 (n.s.)
				A1 : t(19)=1.269	p=2.197e-01 (n.s.)
				S2 : t(19)=0.772	p=4.494e-01 (n.s.)
				A2 : t(19)=0.641	p=5.294e-01 (n.s.)
				S3 : t(19)=-0.108	p=9.149e-01 (n.s.)
				fix` : t(19)=1.404	p=1.763e-01 (n.s.)
				S1` : t(19)=2.393	p=2.718e-02 (*)
V1	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)	fix : t(19)=0.489	p=6.308e-01 (n.s.)
				S1 : t(19)=-0.574	p=5.725e-01 (n.s.)
				A1 : t(19)=-0.452	p=6.565e-01 (n.s.)
				S2 : t(19)=1.674	p=1.105e-01 (n.s.)
				A2 : t(19)=0.221	p=8.272e-01 (n.s.)
				S3 : t(19)=1.458	p=1.612e-01 (n.s.)
				fix` : t(19)=1.275	p=2.176e-01 (n.s.)
				S1` : t(19)=1.088	p=2.901e-01 (n.s.)
HPC	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)	fix : t(19)=0.435	p=6.682e-01 (n.s.)
				S1 : t(19)=1.073	p=2.969e-01 (n.s.)
				A1 : t(19)=0.156	p=8.779e-01 (n.s.)
				S2 : t(19)=0.784	p=4.428e-01 (n.s.)
				A2 : t(19)=1.095	p=2.873e-01 (n.s.)
				S3 : t(19)=1.209	p=2.416e-01 (n.s.)

					fix` : t(19)=1.501	p=1.498e-01 (n.s.)
					S1` : t(19)=-0.183	p=8.568e-01 (n.s.)
vStr 2c, Non-specific goal condition (blue)	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)		fix : t(19)=-0.629	p=5.367e-01 (n.s.)
					S1 : t(19)=1.300	p=2.092e-01 (n.s.)
					A1 : t(19)=2.357	p=2.928e-02 (*)
					S2 : t(19)=2.639	p=1.617e-02 (*)
					A2 : t(19)=-0.036	p=9.720e-01 (n.s.)
					S3 : t(19)=0.588	p=5.633e-01 (n.s.)
					fix` : t(19)=-0.400	p=6.936e-01 (n.s.)
					S1` : t(19)=0.890	p=3.847e-01 (n.s.)
vIPFC 2c, Non-specific goal condition (blue)	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)		fix : t(19)=0.706	p=4.889e-01 (n.s.)
					S1 : t(19)=0.396	p=6.964e-01 (n.s.)
					A1 : t(19)=-0.373	p=7.133e-01 (n.s.)
					S2 : t(19)=0.333	p=7.427e-01 (n.s.)
					A2 : t(19)=1.371	p=1.864e-01 (n.s.)
					S3 : t(19)=1.250	p=2.266e-01 (n.s.)
					fix` : t(19)=0.397	p=6.957e-01 (n.s.)
					S1` : t(19)=1.631	p=1.194e-01 (n.s.)
dlPFC 2c, Non-specific goal condition (blue)	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)		fix : t(19)=0.246	p=8.086e-01 (n.s.)
					S1 : t(19)=-0.153	p=8.797e-01 (n.s.)
					A1 : t(19)=-0.621	p=5.420e-01 (n.s.)
					S2 : t(19)=0.499	p=6.236e-01 (n.s.)
					A2 : t(19)=0.564	p=5.795e-01 (n.s.)
					S3 : t(19)=0.153	p=8.801e-01 (n.s.)
					fix` : t(19)=0.347	p=7.322e-01 (n.s.)
					S1` : t(19)=0.415	p=6.831e-01 (n.s.)
OFC 2c, Non-specific goal condition (blue)	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)		fix : t(19)=0.768	p=4.521e-01 (n.s.)
					S1 : t(19)=0.688	p=4.996e-01 (n.s.)
					A1 : t(19)=0.046	p=9.636e-01 (n.s.)
					S2 : t(19)=1.237	p=2.311e-01 (n.s.)
					A2 : t(19)=0.321	p=7.521e-01 (n.s.)
					S3 : t(19)=0.394	p=6.979e-01 (n.s.)

				fix` : t(19)=0.097	p=9.240e-01 (n.s.)
				S1` : t(19)=1.146	p=2.659e-01 (n.s.)
ACC	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)	fix : t(19)=1.143	p=2.673e-01 (n.s.)
				S1 : t(19)=0.301	p=7.667e-01 (n.s.)
				A1 : t(19)=-1.172	p=2.559e-01 (n.s.)
				S2 : t(19)=0.197	p=8.457e-01 (n.s.)
				A2 : t(19)=0.412	p=6.852e-01 (n.s.)
				S3 : t(19)=1.754	p=9.554e-02 (n.s.)
				fix` : t(19)=0.302	p=7.662e-01 (n.s.)
				S1` : t(19)=0.393	p=6.988e-01 (n.s.)
preSMA	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)	fix : t(19)=0.440	p=6.646e-01 (n.s.)
				S1 : t(19)=1.089	p=2.898e-01 (n.s.)
				A1 : t(19)=0.645	p=5.266e-01 (n.s.)
				S2 : t(19)=1.196	p=2.464e-01 (n.s.)
				A2 : t(19)=1.750	p=9.618e-02 (n.s.)
				S3 : t(19)=0.451	p=6.569e-01 (n.s.)
				fix` : t(19)=0.032	p=9.750e-01 (n.s.)
				S1` : t(19)=0.822	p=4.213e-01 (n.s.)
V1	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)	fix : t(19)=1.332	p=1.988e-01 (n.s.)
				S1 : t(19)=-0.470	p=6.438e-01 (n.s.)
				A1 : t(19)=0.037	p=9.712e-01 (n.s.)
				S2 : t(19)=1.406	p=1.758e-01 (n.s.)
				A2 : t(19)=1.631	p=1.194e-01 (n.s.)
				S3 : t(19)=2.548	p=1.965e-02 (*)
				fix` : t(19)=0.860	p=4.007e-01 (n.s.)
				S1` : t(19)=0.618	p=5.441e-01 (n.s.)
HPC	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)	fix : t(19)=-0.109	p=9.142e-01 (n.s.)
				S1 : t(19)=-0.198	p=8.448e-01 (n.s.)
				A1 : t(19)=1.395	p=1.790e-01 (n.s.)
				S2 : t(19)=1.974	p=6.315e-02 (n.s.)
				A2 : t(19)=1.643	p=1.169e-01 (n.s.)
				S3 : t(19)=2.756	p=1.257e-02 (*)

					fix` : t(19)=1.496 S1` : t(19)=0.889	p=1.512e-01 (n.s.) p=3.850e-01 (n.s.)
vStr	Comparison between event-specific decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)		fix : t(19)=0.505	p=6.197e-01 (n.s.)
					S1 : t(19)=-0.301	p=7.665e-01 (n.s.)
					A1 : t(19)=-0.060	p=9.530e-01 (n.s.)
					S2 : t(19)=1.926	p=6.924e-02 (n.s.)
					A2 : t(19)=1.242	p=2.295e-01 (n.s.)
					S3 : t(19)=0.814	p=4.256e-01 (n.s.)
					fix` : t(19)=1.734	p=9.904e-02 (n.s.)
					S1` : t(19)=-0.030	p=9.764e-01 (n.s.)
2d	Specific goal condition (red) Comparison between event-averaged decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)		vIPFC: t(19)=2.549	p=1.960e-02 (*)
					dIPFC: t(19)=2.412	p=2.612e-02 (*)
					OFC: t(19)=3.749	p=1.359e-03 (**)
					ACC: t(19)=1.521	p=1.448e-01 (n.s.)
					preSMA: t(19)=0.965	p=3.464e-01 (n.s.)
					V1: t(19)=0.748	p=4.633e-01 (n.s.)
					HPC: t(19)=1.513	p=1.468e-01 (n.s.)
					vStr: t(19)=1.617	p=1.224e-01 (n.s.)
Non-specific goal condition (blue)	Comparison between event-averaged decoding accuracy against the chance level (0.5)	20	Paired t-test (two-sided)		vIPFC: t(19)=0.725	p=4.772e-01 (n.s.)
					dIPFC: t(19)=0.175	p=8.629e-01 (n.s.)
					OFC: t(19)=0.546	p=5.913e-01 (n.s.)
					ACC: t(19)=0.394	p=6.983e-01 (n.s.)
					preSMA: t(19)=0.979	p=3.398e-01 (n.s.)
					V1: t(19)=1.290	p=2.126e-01 (n.s.)
					HPC: t(19)=2.093	p=5.002e-02 (n.s.)
					vStr: t(19)=1.185	p=2.505e-01 (n.s.)

Supplementary Table 3 | Statistical details for Fig. 3

Figure panel	Condition	Description	n	Test	Statistic	p-value
3b	vIPFC	Pairwise comparison of event-averaged SDs between four classification types	20	Paired t-test (two-sided)	goal vs unc: $t(19)=1.370$	$p=1.865e-01$ (n.s.)
					goal vs linear: $t(19)=3.322$	$p=3.586e-03$ (**)
					goal vs nonlinear: $t(19)=4.348$	$p=3.464e-04$ (***)
					unc vs linear: $t(19)=-0.122$	$p=9.041e-01$ (n.s.)
					unc vs nonlinear: $t(19)=1.446$	$p=1.644e-01$ (n.s.)
					linear vs nonlinear: $t(19)=3.468$	$p=2.576e-03$ (**)
	dIPFC	Pairwise comparison of event-averaged SDs between four classification types	20	Paired t-test (two-sided)	goal vs unc: $t(19)=2.010$	$p=5.884e-02$ (n.s.)
					goal vs linear: $t(19)=3.679$	$p=1.595e-03$ (**)
					goal vs nonlinear: $t(19)=5.571$	$p=2.258e-05$ (****)
					unc vs linear: $t(19)=-0.806$	$p=4.304e-01$ (n.s.)
					unc vs nonlinear: $t(19)=1.175$	$p=2.544e-01$ (n.s.)
					linear vs nonlinear: $t(19)=4.808$	$p=1.221e-04$ (**)
	OFC	Pairwise comparison of event-averaged SDs between four classification types	20	Paired t-test (two-sided)	goal vs unc: $t(19)=-0.838$	$p=4.125e-01$ (n.s.)
					goal vs linear: $t(19)=1.829$	$p=8.314e-02$ (n.s.)
					goal vs nonlinear: $t(19)=3.453$	$p=2.665e-03$ (**)
					unc vs linear: $t(19)=2.102$	$p=4.909e-02$ (*)
					unc vs nonlinear: $t(19)=2.860$	$p=1.001e-02$ (*)
					linear vs nonlinear: $t(19)=3.250$	$p=4.210e-03$ (**)
	ACC	Pairwise comparison of event-averaged SDs between four classification types	20	Paired t-test (two-sided)	goal vs unc: $t(19)=0.747$	$p=4.641e-01$ (n.s.)
					goal vs linear: $t(19)=1.968$	$p=6.379e-02$ (n.s.)
					goal vs nonlinear: $t(19)=3.342$	$p=3.422e-03$ (**)
					unc vs linear: $t(19)=0.049$	$p=9.612e-01$ (n.s.)
					unc vs nonlinear: $t(19)=1.288$	$p=2.133e-01$ (n.s.)
					linear vs nonlinear: $t(19)=3.465$	$p=2.594e-03$ (**)
	preSMA	Pairwise comparison of event-averaged SDs between four classification types	20	Paired t-test (two-sided)	goal vs unc: $t(19)=1.264$	$p=2.216e-01$ (n.s.)
					goal vs linear: $t(19)=2.330$	$p=3.100e-02$ (*)
					goal vs nonlinear: $t(19)=3.737$	$p=1.395e-03$ (**)
					unc vs linear: $t(19)=-0.592$	$p=5.605e-01$ (n.s.)
					unc vs nonlinear: $t(19)=0.686$	$p=5.011e-01$ (n.s.)
					linear vs nonlinear: $t(19)=3.067$	$p=6.343e-03$ (**)

	V1	Pairwise comparison of event-averaged SDs between four classification types	20	Paired t-test (two-sided)	goal vs unc: $t(19)=2.469$ goal vs linear: $t(19)=4.030$ goal vs nonlinear: $t(19)=5.533$ unc vs linear: $t(19)=-1.561$ unc vs nonlinear: $t(19)=-0.149$ linear vs nonlinear: $t(19)=3.028$	p=2.320e-02 (*) p=7.156e-04 (**) p=2.454e-05 (****) p=1.351e-01 (n.s.) p=8.832e-01 (n.s.) p=6.916e-03 (**)
	HPC	Pairwise comparison of event-averaged SDs between four classification types	20	Paired t-test (two-sided)	goal vs unc: $t(19)=-1.021$ goal vs linear: $t(19)=0.277$ goal vs nonlinear: $t(19)=1.618$ unc vs linear: $t(19)=1.537$ unc vs nonlinear: $t(19)=1.616$ linear vs nonlinear: $t(19)=1.490$	p=3.199e-01 (n.s.) p=7.850e-01 (n.s.) p=1.221e-01 (n.s.) p=1.408e-01 (n.s.) p=1.226e-01 (n.s.) p=1.525e-01 (n.s.)
	vStr	Pairwise comparison of event-averaged SDs between four classification types	20	Paired t-test (two-sided)	goal vs unc: $t(19)=-1.201$ goal vs linear: $t(19)=-0.655$ goal vs nonlinear: $t(19)=-0.243$ unc vs linear: $t(19)=1.428$ unc vs nonlinear: $t(19)=1.321$ linear vs nonlinear: $t(19)=0.507$	p=2.446e-01 (n.s.) p=5.204e-01 (n.s.) p=8.107e-01 (n.s.) p=1.696e-01 (n.s.) p=2.023e-01 (n.s.) p=6.181e-01 (n.s.)
3d, 1st row (vs Behavioral flexibility)	Goal SD	Correlation between neural goal SD and average choice versatility across participants, corrected for the number of ROIs (Benjamini-Hochberg procedure, q=0.05)	20	Pearson's correlation	vIPFC: $r=0.583$, $t(18)=3.045$ dIPFC: $r=0.641$, $t(18)=3.545$ OFC: $r=0.443$, $t(18)=2.094$ ACC: $r=0.453$, $t(18)=2.154$ preSMA: $r=0.447$, $t(18)=2.119$ V1: $r=0.305$, $t(18)=1.360$ HPC: $r=0.002$, $t(18)=0.009$ vStr: $r=0.351$, $t(18)=1.588$	p=2.785e-02 (*) p=1.852e-02 (*) p=8.102e-02 (n.s.) p=8.102e-02 (n.s.) p=8.102e-02 (n.s.) p=2.178e-01 (n.s.) p=9.930e-01 (n.s.) p=1.730e-01 (n.s.)
	Uncertainty SD	Correlation between neural goal SD and average choice versatility across participants, corrected for the number of ROIs (Benjamini-Hochberg procedure, q=0.05)	20	Pearson's correlation	vIPFC: $r=-0.150$, $t(18)=-0.646$ dIPFC: $r=-0.050$, $t(18)=-0.212$ OFC: $r=0.226$, $t(18)=0.985$ ACC: $r=-0.055$, $t(18)=-0.234$ preSMA: $r=-0.226$, $t(18)=-0.982$ V1: $r=0.047$, $t(18)=0.198$	p=8.454e-01 (n.s.) p=8.454e-01 (n.s.) p=8.454e-01 (n.s.) p=8.454e-01 (n.s.) p=8.454e-01 (n.s.) p=8.454e-01 (n.s.)

					HPC: r=0.396, t(18)=1.830 vStr: r=0.056, t(18)=0.237	p=6.714e-01 (n.s.) p=8.454e-01 (n.s.)
Linear SD	Correlation between neural goal SD and average choice versatility across participants, corrected for the number of ROIs (Benjamini-Hochberg procedure, q=0.05)	20	Pearson's correlation	vIPFC: r=0.370, t(18)=1.692 dIPFC: r=0.465, t(18)=2.225 OFC: r=0.441, t(18)=2.083 ACC: r=0.236, t(18)=1.030 preSMA: r=0.189, t(18)=0.818 V1: r=0.210, t(18)=0.911 HPC: r=0.281, t(18)=1.240 vStr: r=0.387, t(18)=1.780	p=2.157e-01 (n.s.) p=2.072e-01 (n.s.) p=2.072e-01 (n.s.) p=4.224e-01 (n.s.) p=4.243e-01 (n.s.) p=4.243e-01 (n.s.) p=3.695e-01 (n.s.) p=2.157e-01 (n.s.)	
				vIPFC: r=0.399, t(18)=1.844 dIPFC: r=0.537, t(18)=2.703 OFC: r=0.311, t(18)=1.386 ACC: r=-0.049, t(18)=-0.206 preSMA: r=0.291, t(18)=1.292 V1: r=0.277, t(18)=1.223 HPC: r=0.059, t(18)=0.252 vStr: r=0.556, t(18)=2.839	p=2.177e-01 (n.s.) p=5.825e-02 (n.s.) p=3.164e-01 (n.s.) p=8.390e-01 (n.s.) p=3.164e-01 (n.s.) p=3.164e-01 (n.s.) p=8.390e-01 (n.s.) p=5.825e-02 (n.s.)	
				vIPFC: r=0.554, t(18)=2.823 dIPFC: r=0.577, t(18)=2.994 OFC: r=0.422, t(18)=1.977 ACC: r=0.262, t(18)=1.152 preSMA: r=0.553, t(18)=2.813 V1: r=0.282, t(18)=1.248 HPC: r=0.191, t(18)=0.825 vStr: r=0.197, t(18)=0.852	p=3.067e-02 (*) p=3.067e-02 (*) p=1.272e-01 (n.s.) p=3.523e-01 (n.s.) p=3.067e-02 (*) p=3.523e-01 (n.s.) p=4.202e-01 (n.s.) p=4.202e-01 (n.s.)	
				vIPFC: r=-0.196, t(18)=-0.849 dIPFC: r=-0.328, t(18)=-1.471 OFC: r=0.154, t(18)=0.661 ACC: r=0.017, t(18)=0.070 preSMA: r=-0.211, t(18)=-0.914 V1: r=-0.002, t(18)=-0.007	p=8.138e-01 (n.s.) p=8.138e-01 (n.s.) p=8.276e-01 (n.s.) p=9.942e-01 (n.s.) p=8.138e-01 (n.s.) p=9.942e-01 (n.s.)	
3d, 2nd row (vs Behavioral stability)	Correlation between neural goal SD and average choice consistency across participants, corrected for the number of ROIs (Benjamini-Hochberg procedure, q=0.05)	20	Pearson's correlation	vIPFC: r=0.554, t(18)=2.823 dIPFC: r=0.577, t(18)=2.994 OFC: r=0.422, t(18)=1.977 ACC: r=0.262, t(18)=1.152 preSMA: r=0.553, t(18)=2.813 V1: r=0.282, t(18)=1.248 HPC: r=0.191, t(18)=0.825 vStr: r=0.197, t(18)=0.852	p=3.067e-02 (*) p=3.067e-02 (*) p=1.272e-01 (n.s.) p=3.523e-01 (n.s.) p=3.067e-02 (*) p=3.523e-01 (n.s.) p=4.202e-01 (n.s.) p=4.202e-01 (n.s.)	
				vIPFC: r=0.554, t(18)=2.823 dIPFC: r=0.577, t(18)=2.994 OFC: r=0.422, t(18)=1.977 ACC: r=0.262, t(18)=1.152 preSMA: r=0.553, t(18)=2.813 V1: r=0.282, t(18)=1.248 HPC: r=0.191, t(18)=0.825 vStr: r=0.197, t(18)=0.852	p=3.067e-02 (*) p=3.067e-02 (*) p=1.272e-01 (n.s.) p=3.523e-01 (n.s.) p=3.067e-02 (*) p=3.523e-01 (n.s.) p=4.202e-01 (n.s.) p=4.202e-01 (n.s.)	
				vIPFC: r=0.554, t(18)=2.823 dIPFC: r=0.577, t(18)=2.994 OFC: r=0.422, t(18)=1.977 ACC: r=0.262, t(18)=1.152 preSMA: r=0.553, t(18)=2.813 V1: r=0.282, t(18)=1.248 HPC: r=0.191, t(18)=0.825 vStr: r=0.197, t(18)=0.852	p=3.067e-02 (*) p=3.067e-02 (*) p=1.272e-01 (n.s.) p=3.523e-01 (n.s.) p=3.067e-02 (*) p=3.523e-01 (n.s.) p=4.202e-01 (n.s.) p=4.202e-01 (n.s.)	
				vIPFC: r=0.554, t(18)=2.823 dIPFC: r=0.577, t(18)=2.994 OFC: r=0.422, t(18)=1.977 ACC: r=0.262, t(18)=1.152 preSMA: r=0.553, t(18)=2.813 V1: r=0.282, t(18)=1.248 HPC: r=0.191, t(18)=0.825 vStr: r=0.197, t(18)=0.852	p=3.067e-02 (*) p=3.067e-02 (*) p=1.272e-01 (n.s.) p=3.523e-01 (n.s.) p=3.067e-02 (*) p=3.523e-01 (n.s.) p=4.202e-01 (n.s.) p=4.202e-01 (n.s.)	
				vIPFC: r=0.554, t(18)=2.823 dIPFC: r=0.577, t(18)=2.994 OFC: r=0.422, t(18)=1.977 ACC: r=0.262, t(18)=1.152 preSMA: r=0.553, t(18)=2.813 V1: r=0.282, t(18)=1.248 HPC: r=0.191, t(18)=0.825 vStr: r=0.197, t(18)=0.852	p=3.067e-02 (*) p=3.067e-02 (*) p=1.272e-01 (n.s.) p=3.523e-01 (n.s.) p=3.067e-02 (*) p=3.523e-01 (n.s.) p=4.202e-01 (n.s.) p=4.202e-01 (n.s.)	
				vIPFC: r=0.554, t(18)=2.823 dIPFC: r=0.577, t(18)=2.994 OFC: r=0.422, t(18)=1.977 ACC: r=0.262, t(18)=1.152 preSMA: r=0.553, t(18)=2.813 V1: r=0.282, t(18)=1.248 HPC: r=0.191, t(18)=0.825 vStr: r=0.197, t(18)=0.852	p=3.067e-02 (*) p=3.067e-02 (*) p=1.272e-01 (n.s.) p=3.523e-01 (n.s.) p=3.067e-02 (*) p=3.523e-01 (n.s.) p=4.202e-01 (n.s.) p=4.202e-01 (n.s.)	
				vIPFC: r=0.554, t(18)=2.823 dIPFC: r=0.577, t(18)=2.994 OFC: r=0.422, t(18)=1.977 ACC: r=0.262, t(18)=1.152 preSMA: r=0.553, t(18)=2.813 V1: r=0.282, t(18)=1.248 HPC: r=0.191, t(18)=0.825 vStr: r=0.197, t(18)=0.852	p=3.067e-02 (*) p=3.067e-02 (*) p=1.272e-01 (n.s.) p=3.523e-01 (n.s.) p=3.067e-02 (*) p=3.523e-01 (n.s.) p=4.202e-01 (n.s.) p=4.202e-01 (n.s.)	
				vIPFC: r=0.554, t(18)=2.823 dIPFC: r=0.577, t(18)=2.994 OFC: r=0.422, t(18)=1.977 ACC: r=0.262, t(18)=1.152 preSMA: r=0.553, t(18)=2.813 V1: r=0.282, t(18)=1.248 HPC: r=0.191, t(18)=0.825 vStr: r=0.197, t(18)=0.852	p=3.067e-02 (*) p=3.067e-02 (*) p=1.272e-01 (n.s.) p=3.523e-01 (n.s.) p=3.067e-02 (*) p=3.523e-01 (n.s.) p=4.202e-01 (n.s.) p=4.202e-01 (n.s.)	

					HPC: r=0.251, t(18)=1.102 vStr: r=-0.062, t(18)=-0.264	p=8.138e-01 (n.s.) p=9.942e-01 (n.s.)
Linear SD	Correlation between neural goal SD and average choice consistency across participants, corrected for the number of ROIs (Benjamini-Hochberg procedure, q=0.05)	20	Pearson's correlation	vIPFC: r=0.319, t(18)=1.427 dIPFC: r=0.238, t(18)=1.041 OFC: r=0.341, t(18)=1.539 ACC: r=0.160, t(18)=0.689 preSMA: r=0.251, t(18)=1.101 V1: r=0.156, t(18)=0.671 HPC: r=0.278, t(18)=1.229 vStr: r=0.269, t(18)=1.183	p=4.153e-01 (n.s.) p=4.153e-01 (n.s.) p=4.153e-01 (n.s.) p=5.104e-01 (n.s.) p=4.153e-01 (n.s.) p=5.104e-01 (n.s.) p=4.153e-01 (n.s.) p=4.153e-01 (n.s.)	
				vIPFC: r=0.210, t(18)=0.910 dIPFC: r=0.305, t(18)=1.359 OFC: r=0.058, t(18)=0.247 ACC: r=-0.114, t(18)=-0.487 preSMA: r=0.224, t(18)=0.975 V1: r=0.041, t(18)=0.173 HPC: r=0.063, t(18)=0.269 vStr: r=0.470, t(18)=2.259	p=7.494e-01 (n.s.) p=7.494e-01 (n.s.) p=8.644e-01 (n.s.) p=8.644e-01 (n.s.) p=7.494e-01 (n.s.) p=8.644e-01 (n.s.) p=8.644e-01 (n.s.) p=2.922e-01 (n.s.)	
				vIPFC: r=0.679, t(18)=3.929 dIPFC: r=0.768, t(18)=5.084 OFC: r=0.549, t(18)=2.784 ACC: r=0.542, t(18)=2.740 preSMA: r=0.656, t(18)=3.690 V1: r=0.414, t(18)=1.930 HPC: r=0.257, t(18)=1.126 vStr: r=0.292, t(18)=1.293	p=3.933e-03 (**) p=6.201e-04 (***) p=2.154e-02 (*) p=2.154e-02 (*) p=4.471e-03 (**) p=9.268e-02 (n.s.) p=2.748e-01 (n.s.) p=2.427e-01 (n.s.)	
				vIPFC: r=-0.143, t(18)=-0.611 dIPFC: r=-0.095, t(18)=-0.404 OFC: r=0.279, t(18)=1.234 ACC: r=0.211, t(18)=0.915 preSMA: r=-0.100, t(18)=-0.428 V1: r=0.174, t(18)=0.752	p=7.889e-01 (n.s.) p=7.889e-01 (n.s.) p=7.889e-01 (n.s.) p=7.889e-01 (n.s.) p=7.889e-01 (n.s.) p=7.889e-01 (n.s.)	
				vIPFC: r=0.319, t(18)=1.427 dIPFC: r=0.238, t(18)=1.041 OFC: r=0.341, t(18)=1.539 ACC: r=0.160, t(18)=0.689 preSMA: r=0.251, t(18)=1.101 V1: r=0.156, t(18)=0.671 HPC: r=0.278, t(18)=1.229 vStr: r=0.269, t(18)=1.183	p=4.153e-01 (n.s.) p=4.153e-01 (n.s.) p=4.153e-01 (n.s.) p=5.104e-01 (n.s.) p=4.153e-01 (n.s.) p=5.104e-01 (n.s.) p=4.153e-01 (n.s.) p=4.153e-01 (n.s.)	
				vIPFC: r=0.210, t(18)=0.910 dIPFC: r=0.305, t(18)=1.359 OFC: r=0.058, t(18)=0.247 ACC: r=-0.114, t(18)=-0.487 preSMA: r=0.224, t(18)=0.975 V1: r=0.041, t(18)=0.173 HPC: r=0.063, t(18)=0.269 vStr: r=0.470, t(18)=2.259	p=7.494e-01 (n.s.) p=7.494e-01 (n.s.) p=8.644e-01 (n.s.) p=8.644e-01 (n.s.) p=7.494e-01 (n.s.) p=8.644e-01 (n.s.) p=8.644e-01 (n.s.) p=2.922e-01 (n.s.)	
				vIPFC: r=0.679, t(18)=3.929 dIPFC: r=0.768, t(18)=5.084 OFC: r=0.549, t(18)=2.784 ACC: r=0.542, t(18)=2.740 preSMA: r=0.656, t(18)=3.690 V1: r=0.414, t(18)=1.930 HPC: r=0.257, t(18)=1.126 vStr: r=0.292, t(18)=1.293	p=3.933e-03 (**) p=6.201e-04 (***) p=2.154e-02 (*) p=2.154e-02 (*) p=4.471e-03 (**) p=9.268e-02 (n.s.) p=2.748e-01 (n.s.) p=2.427e-01 (n.s.)	
				vIPFC: r=-0.143, t(18)=-0.611 dIPFC: r=-0.095, t(18)=-0.404 OFC: r=0.279, t(18)=1.234 ACC: r=0.211, t(18)=0.915 preSMA: r=-0.100, t(18)=-0.428 V1: r=0.174, t(18)=0.752	p=7.889e-01 (n.s.) p=7.889e-01 (n.s.) p=7.889e-01 (n.s.) p=7.889e-01 (n.s.) p=7.889e-01 (n.s.) p=7.889e-01 (n.s.)	
3d, 3rd row (vs Behavioral performance)	Goal SD	Correlation between neural goal SD and average choice optimality across participants, corrected for the number of ROIs (Benjamini-Hochberg procedure, q=0.05)	20	Pearson's correlation	vIPFC: r=0.679, t(18)=3.929 dIPFC: r=0.768, t(18)=5.084 OFC: r=0.549, t(18)=2.784 ACC: r=0.542, t(18)=2.740 preSMA: r=0.656, t(18)=3.690 V1: r=0.414, t(18)=1.930 HPC: r=0.257, t(18)=1.126 vStr: r=0.292, t(18)=1.293	p=3.933e-03 (**) p=6.201e-04 (***) p=2.154e-02 (*) p=2.154e-02 (*) p=4.471e-03 (**) p=9.268e-02 (n.s.) p=2.748e-01 (n.s.) p=2.427e-01 (n.s.)
Uncertainty SD	Correlation between neural goal SD and average choice optimality across participants, corrected for the number of ROIs (Benjamini-Hochberg procedure, q=0.05)	20	Pearson's correlation	vIPFC: r=-0.143, t(18)=-0.611 dIPFC: r=-0.095, t(18)=-0.404 OFC: r=0.279, t(18)=1.234 ACC: r=0.211, t(18)=0.915 preSMA: r=-0.100, t(18)=-0.428 V1: r=0.174, t(18)=0.752	p=7.889e-01 (n.s.) p=7.889e-01 (n.s.) p=7.889e-01 (n.s.) p=7.889e-01 (n.s.) p=7.889e-01 (n.s.) p=7.889e-01 (n.s.)	

					HPC: $r=0.365$, $t(18)=1.665$	$p=7.889e-01$ (n.s.)
					vStr: $r=-0.064$, $t(18)=-0.272$	$p=7.889e-01$ (n.s.)
Linear SD	Correlation between neural goal SD and average choice optimality across participants, corrected for the number of ROIs (Benjamini-Hochberg procedure, $q=0.05$)	20	Pearson's correlation	vIPFC: $r=0.485$, $t(18)=2.354$	$p=8.039e-02$ (n.s.)	
				dIPFC: $r=0.491$, $t(18)=2.389$	$p=8.039e-02$ (n.s.)	
				OFC: $r=0.552$, $t(18)=2.810$	$p=8.039e-02$ (n.s.)	
				ACC: $r=0.445$, $t(18)=2.110$	$p=9.830e-02$ (n.s.)	
				preSMA: $r=0.364$, $t(18)=1.660$	$p=1.465e-01$ (n.s.)	
				V1: $r=0.352$, $t(18)=1.595$	$p=1.465e-01$ (n.s.)	
				HPC: $r=0.374$, $t(18)=1.709$	$p=1.465e-01$ (n.s.)	
				vStr: $r=0.299$, $t(18)=1.329$	$p=2.003e-01$ (n.s.)	
Nonlinear SD	Correlation between neural goal SD and average choice optimality across participants, corrected for the number of ROIs (Benjamini-Hochberg procedure, $q=0.05$)	20	Pearson's correlation	vIPFC: $r=0.422$, $t(18)=1.973$	$p=1.708e-01$ (n.s.)	
				dIPFC: $r=0.448$, $t(18)=2.125$	$p=1.708e-01$ (n.s.)	
				OFC: $r=0.296$, $t(18)=1.316$	$p=4.094e-01$ (n.s.)	
				ACC: $r=-0.064$, $t(18)=-0.272$	$p=7.884e-01$ (n.s.)	
				preSMA: $r=0.259$, $t(18)=1.139$	$p=4.315e-01$ (n.s.)	
				V1: $r=0.158$, $t(18)=0.680$	$p=6.737e-01$ (n.s.)	
				HPC: $r=0.066$, $t(18)=0.280$	$p=7.884e-01$ (n.s.)	
				vStr: $r=0.464$, $t(18)=2.220$	$p=1.708e-01$ (n.s.)	

Supplementary Table 4 | Statistical details for Fig. 4

Figure panel	Condition	Description	n	Test	Statistic	p-value
4a	vIPFC	Pairwise comparison of event-averaged neural measures	20	Paired t-test (two-sided)	CCGP vs low uncert. SD: $t(19)=-0.337$	$p=7.400e-01$ (n.s.)
					CCGP vs high uncert. SD: $t(19)=-0.827$	$p=4.186e-01$ (n.s.)
					low uncert. SD vs high uncert. SD: $t(19)=-0.180$	$p=8.594e-01$ (n.s.)
	dIPFC	Pairwise comparison of event-averaged neural measures	20	Paired t-test (two-sided)	CCGP vs low uncert. SD: $t(19)=-0.072$	$p=9.434e-01$ (n.s.)
					CCGP vs high uncert. SD: $t(19)=-0.372$	$p=7.142e-01$ (n.s.)
					low uncert. SD vs high uncert. SD: $t(19)=-0.216$	$p=8.316e-01$ (n.s.)
	OFC	Pairwise comparison of event-averaged neural measures	20	Paired t-test (two-sided)	CCGP vs low uncert. SD: $t(19)=0.109$	$p=9.146e-01$ (n.s.)
					CCGP vs high uncert. SD: $t(19)=-0.395$	$p=6.970e-01$ (n.s.)
					low uncert. SD vs high uncert. SD: $t(19)=-0.513$	$p=6.135e-01$ (n.s.)
	ACC	Pairwise comparison of event-averaged neural measures	20	Paired t-test (two-sided)	CCGP vs low uncert. SD: $t(19)=-0.038$	$p=9.698e-01$ (n.s.)
					CCGP vs high uncert. SD: $t(19)=0.405$	$p=6.902e-01$ (n.s.)
					low uncert. SD vs high uncert. SD: $t(19)=0.423$	$p=6.768e-01$ (n.s.)
	preSMA	Pairwise comparison of event-averaged neural measures	20	Paired t-test (two-sided)	CCGP vs low uncert. SD: $t(19)=0.315$	$p=7.563e-01$ (n.s.)
					CCGP vs high uncert. SD: $t(19)=0.873$	$p=3.937e-01$ (n.s.)
					low uncert. SD vs high uncert. SD: $t(19)=0.395$	$p=6.976e-01$ (n.s.)
	V1	Pairwise comparison of event-averaged neural measures	20	Paired t-test (two-sided)	CCGP vs low uncert. SD: $t(19)=-1.422$	$p=1.713e-01$ (n.s.)
					CCGP vs high uncert. SD: $t(19)=-0.549$	$p=5.893e-01$ (n.s.)
					low uncert. SD vs high uncert. SD: $t(19)=0.857$	$p=4.021e-01$ (n.s.)
	HPC	Pairwise comparison of event-averaged neural measures	20	Paired t-test (two-sided)	CCGP vs low uncert. SD: $t(19)=0.069$	$p=9.457e-01$ (n.s.)
					CCGP vs high uncert. SD: $t(19)=1.778$	$p=9.142e-02$ (n.s.)
					low uncert. SD vs high uncert. SD: $t(19)=1.234$	$p=2.322e-01$ (n.s.)
	vStr	Pairwise comparison of event-averaged neural measures	20	Paired t-test (two-sided)	CCGP vs low uncert. SD: $t(19)=-0.973$	$p=3.428e-01$ (n.s.)
					CCGP vs high uncert. SD: $t(19)=-0.382$	$p=7.064e-01$ (n.s.)
					low uncert. SD vs high uncert. SD: $t(19)=0.439$	$p=6.654e-01$ (n.s.)
4b	Left (CCGP vs behavioral flexibility)	Correlation between neural goal CCGP and average choice versatility across participants, corrected for the number of ROIs (Benjamini-Hochberg procedure, $q=0.05$)	20	Pearson's correlation	vIPFC: $r=0.576$, $t(18)=2.990$	$p=4.703e-02$ (*)
					dIPFC: $r=0.551$, $t(18)=2.803$	$p=4.703e-02$ (*)
					OFC: $r=0.293$, $t(18)=1.298$	$p=3.371e-01$ (n.s.)
					ACC: $r=0.472$, $t(18)=2.269$	$p=9.539e-02$ (n.s.)
					preSMA: $r=0.387$, $t(18)=1.783$	$p=1.831e-01$ (n.s.)
					V1: $r=0.023$, $t(18)=0.097$	$p=9.240e-01$ (n.s.)
					HPC: $r=0.162$, $t(18)=0.698$	$p=6.585e-01$ (n.s.)

				vStr: r=0.034, t(18)=0.143	p=9.240e-01 (n.s.)
Middle (CCGP vs behavioral stability)	Correlation between neural goal CCGP and average choice consistency across participants, corrected for the number of ROIs (Benjamini-Hochberg procedure, q=0.05)	20	Pearson's correlation	vIPFC: r=0.559, t(18)=2.859	p=5.364e-02 (n.s.)
				dIPFC: r=0.543, t(18)=2.742	p=5.364e-02 (n.s.)
				OFC: r=0.468, t(18)=2.245	p=7.514e-02 (n.s.)
				ACC: r=0.341, t(18)=1.538	p=2.264e-01 (n.s.)
				preSMA: r=0.511, t(18)=2.525	p=5.642e-02 (n.s.)
				V1: r=0.184, t(18)=0.796	p=4.989e-01 (n.s.)
				HPC: r=0.221, t(18)=0.959	p=4.668e-01 (n.s.)
				vStr: r=-0.073, t(18)=-0.311	p=7.592e-01 (n.s.)
Right (CCGP vs behavioral performance)	Correlation between neural goal CCGP and average choice optimality across participants, corrected for the number of ROIs (Benjamini-Hochberg procedure, q=0.05)	20	Pearson's correlation	vIPFC: r=0.628, t(18)=3.427	p=8.024e-03 (**)
				dIPFC: r=0.713, t(18)=4.314	p=3.342e-03 (**)
				OFC: r=0.457, t(18)=2.181	p=6.832e-02 (n.s.)
				ACC: r=0.564, t(18)=2.901	p=1.903e-02 (*)
				preSMA: r=0.641, t(18)=3.541	p=8.024e-03 (**)
				V1: r=0.235, t(18)=1.024	p=3.651e-01 (n.s.)
				HPC: r=0.281, t(18)=1.244	p=3.057e-01 (n.s.)
				vStr: r=-0.003, t(18)=-0.012	p=9.907e-01 (n.s.)

Please wait...

If this message is not eventually replaced by the proper contents of the document, your PDF viewer may not be able to display this type of document.

You can upgrade to the latest version of Adobe Reader for Windows®, Mac, or Linux® by visiting http://www.adobe.com/go/reader_download.

For more assistance with Adobe Reader visit <http://www.adobe.com/go/acrreader>.

Windows is either a registered trademark or a trademark of Microsoft Corporation in the United States and/or other countries. Mac is a trademark of Apple Inc., registered in the United States and other countries. Linux is the registered trademark of Linus Torvalds in the U.S. and other countries.