

FERAL: A Video-Understanding System for Direct Video-to-Behavior Mapping

Peter Skovorodnikov^{1,†,*}, Janet Zhao², Friederike Buck³, Tomas Kay², Dominic D. Frank², Benjamin Koger^{4,5}, Blair R. Costelloe^{6,7,8}, Iain D. Couzin^{6,7,8}, Jacopo Razzauti^{9,10,11,†,*}

¹Data Science Platform, The Rockefeller University, New York, NY, USA

²Laboratory of Social Evolution and Behavior, The Rockefeller University, New York, NY, USA

³Lulu and Anthony Wang Laboratory of Neural Circuits and Behavior, The Rockefeller University, New York, NY, USA

⁴School of Computing, University of Wyoming, Laramie, WY, USA

⁵Department of Zoology and Physiology, University of Wyoming, Laramie, WY, USA

⁶Department of Collective Behaviour, Max Planck Institute of Animal Behavior, Konstanz, Germany

⁷Department of Biology, University of Konstanz, Konstanz, Germany

⁸Centre for the Advanced Study of Collective Behaviour, University of Konstanz, Germany

⁹Laboratory of Neurogenetics and Behavior, The Rockefeller University, New York, NY, USA

¹⁰The Price Family Center for the Social Brain, The Rockefeller University, New York, NY, USA

¹¹Howard Hughes Medical Institute, New York, NY, USA

Correspondence: peter.skovorodnikov@gmail.com, jacopo.razza@gmail.com

[†]These authors contributed equally.

1

Abstract

2

Animal behavior unfolds continuously in time, yet quantitative analyses often require segmenting it into discrete, interpretable states. Although manual annotation can achieve this, it remains slow, subjective, and difficult to scale. Most automated pipelines use tracked body parts to infer actions, but are limited by tracking quality, and discard much of the visual information contained in raw videos. Here we present FERAL (Feature Extraction for Recognition of Animal Locomotion), a supervised video-understanding toolkit that bridges this gap by mapping raw video directly to frame-level behavioral labels, bypassing the need for pose estimation. Across benchmarks, FERAL outperforms state-of-the-art pose- and video-based baselines: on a benchmarking dataset of mouse social interaction, it surpasses Google’s Videoprism using just a quarter of the training data. FERAL generalizes across species, recording conditions, and

12 levels of behavioral organization: from single-animal locomotion to complex social interactions
13 and emergent collective dynamics. Released as a user-friendly, open-source package, FERAL
14 overcomes the challenges of traditional approaches, integrates easily with existing analysis
15 pipelines, and can be deployed locally or on cloud servers with a few clicks. By mapping
16 raw video directly to annotated behavior, FERAL lowers the barrier to scalable, cross-species
17 behavioral quantification and broadens the range of behavioral analyses possible in both the lab
18 and the wild.

19 **Introduction**

20 Quantifying animal behavior remains a central challenge in biology, linking neural activity, evolution,
21 and ecology through the study of discrete, observable actions [1–4]. Since the earliest days of
22 ethology, researchers have identified these actions through naturalistic observation and manual
23 annotation [5–9], establishing the foundation of modern behavioral science [3, 10–12]. Such
24 analyses have been essential across disciplines, informing research in ecology [13], neuroscience
25 [14], disease modeling, and drug discovery.

26 Manual annotation remains the most versatile method for quantifying behavior, adapting to scene
27 complexities that range from controlled arenas to field recordings. However, it is labor-intensive,
28 subjective, and the main bottleneck in scaling behavioral research [15–17]. Advances in computer
29 vision have partially mitigated this limitation through markerless tracking tools such as DeepLabCut
30 [18], SLEAP [19], and LightningPose [20], which estimate the positions of animals and their body
31 parts over time [21, 22]. These tools have transformed behavioral analysis in laboratory settings,
32 but still require tightly controlled imaging conditions. In more naturalistic or noisy environments,
33 they often fail, forcing researchers to simplify behavioral assays (e.g., by tethering animals or
34 designing task-specific arenas), rely on manual annotation, or to abandon complex analyses altogether
35 (**Figure 1a-b**). Importantly, these methods answer where animals (and their body parts) are, but not
36 what actions they are performing.

37 Automatically determining the actions animals are engaged in remains a greater challenge.

38 Current computational pipelines for behavioral segmentation depend almost entirely on pose-based
39 outputs, transforming keypoint trajectories into discrete behavioral states [2, 23–28]. Analyses based
40 on these intermediate, skeletonized representations are limited by tracking quality and feasibility,
41 and discard rich contextual information that is often essential for accurate behavioral interpretation.
42 As the level of behavioral organization increases, such as in social or collective dynamics, these
43 pipelines become increasingly complex, requiring multi-animal tracking and extensive post-tracking
44 curation (**Figure 1a**). In many cases, the visual cues that define a behavior are well established, and
45 researchers primarily need a scalable, automated way to recognize them across large datasets.

46 To bridge this gap, we developed FERAL (Feature Extraction for Recognition of Animal
47 Locomotion), a supervised video-understanding system that learns behavior directly from raw video
48 frames, bypassing the need for pose trajectories while preserving full temporal and visual information
49 (**Figure 1b**). This approach is particularly advantageous when pose estimation is unreliable or
50 unnecessary for the analysis at hand. FERAL can operate either alongside or independently of
51 existing pose-based frameworks, providing an end-to-end solution for context-rich behavioral
52 segmentation.

53 Leveraging a pretrained video foundation model [29], FERAL integrates motion and visual cues
54 within a unified architecture. Trained directly on videos aligned with expert annotations, it detects
55 discrete actions and generates interpretable behavioral sequences (i.e., ethograms) that reliably
56 capture discrete animal behavior.

57 On benchmarking datasets, FERAL consistently outperforms state-of-the-art pose- and video-
58 based baselines. It generalizes robustly across seven datasets spanning diverse species, recording
59 modalities (from controlled laboratory conditions to field and aerial footages), and levels of
60 behavioral organization: from single-animal locomotion in *C. elegans* to social interactions and
61 emergent collective dynamics in rodents, ants, and primates. These results demonstrate that FERAL
62 maintains high performance even in scenarios where traditional pipelines typically struggle. For
63 instance, FERAL makes analysis of behaviors involving extensive interaction of animals either with
64 each other or with their environments much more feasible.

65 Designed for ease of use, FERAL unifies preprocessing, training, and inference within a modular,
66 user-friendly workflow requiring no coding or machine learning expertise. With only a few lines
67 of Python it can be run locally or on cloud platforms such as Google Colab, providing a scalable,
68 context-aware foundation for modern behavioral science.

69 **Results**

70 **Direct video-to-behavior mapping with FERAL**

71 FERAL provides a user-friendly, end-to-end workflow that takes raw videos and behavioral
72 annotations as inputs, fine-tunes an open-source video-understanding model, and outputs segmented
73 behavioral sequences. To ensure reproducibility across laboratories and recording setups, the
74 pipeline standardizes both inputs (videos and labels) into a unified format.

75 The first stage is *video preprocessing*. Because video-understanding models operate on relatively
76 low-resolution inputs (e.g., 512×512 pixels), input videos must be resized and re-encoded into
77 seekable formats to enable efficient frame sampling [29, 30]. FERAL includes a cross-platform
78 function that automates this process, producing standardized inputs regardless of the original
79 recording format (**Figure 1c**).

80 The second stage, *label preparation*, converts behavioral annotations from diverse software
81 (e.g., BORIS [31], EthoVision [32]) into a consistent schema [31, 33]. The annotations provided
82 by the user, often timestamp-based, are converted into JSON files that map every video frame to
83 a categorical behavioral label (**Figure 1d**). In practice, users need only a folder of videos and a
84 corresponding label file to initiate training and inference.

85 Because full-length recordings cannot be processed in a single pass, FERAL divides each video
86 into overlapping temporal chunks that are independently processed and corresponding predictions
87 are subsequently ensembled (**Figure 1e**). This design minimizes misclassification at behavioral
88 boundaries and maintains temporal continuity by ensuring that each short behavior is fully captured
89 within at least one chunk.

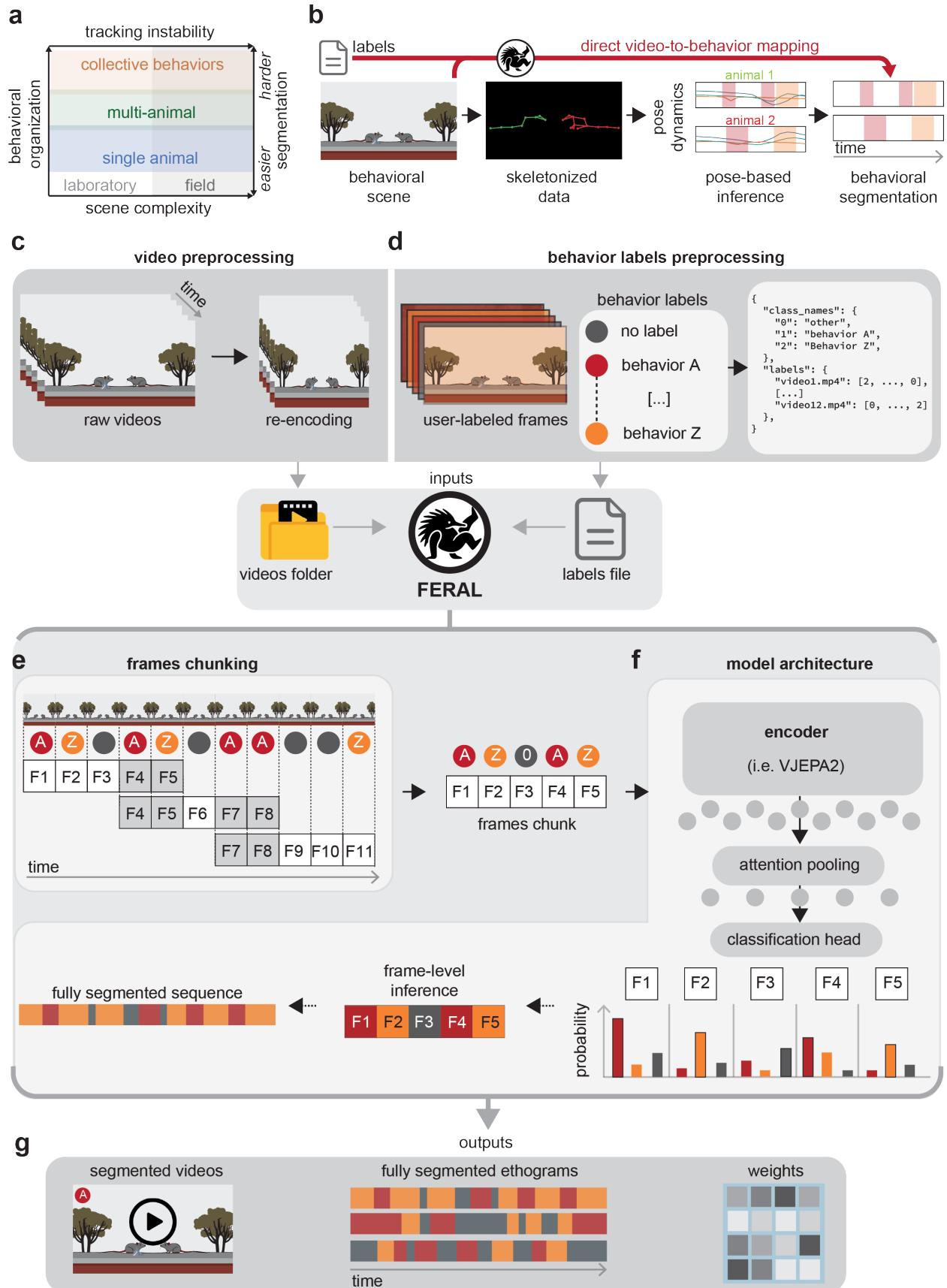


Figure 1: **Overview of the FERAL workflow** **(a)** Conceptual space of behavioral analysis: increasing scene complexity (left to right) and level of organization (bottom to top) challenge tracking and pose-based pipelines for behavior segmentation. **(b)** Direct video-to-behavior mapping: FERAL learns from raw scenes and expert labels, bypassing skeletonized pose to produce frame-level segmentations. **(c)** *Video preprocessing*: raw videos are resized and re-encoded into standardized, seekable inputs. **(d)** *Behavior labels preprocessing*: user annotations are converted to a unified frame-aligned JSON schema (class names and per-frame label arrays). **(e)** *Frames chunking*: videos are divided into overlapping temporal segments. **(f)** *Model training*: chunks are embedded by a pretrained video encoder (i.e., V-JEPA2); attention pooling aggregates spatiotemporal features followed by a lightweight head. Overlapping predictions are ensembled to produce stable frame-level probabilities. **(g)** *Outputs*: segmented videos, ethograms, and model weights for reuse or further fine-tuning.

90 After evaluating several potential backbones, we adopted V-JEPA2 [29], a recent foundation
91 model introduced by Meta FAIR (see Methods). Trained self-supervised on over one million hours
92 of unlabeled video, V-JEPA2 learns rich spatiotemporal embeddings through a masked-prediction
93 objective, producing a robust representation of input videos. Within FERAL, we fine-tune only a
94 subset of its transformer layers using comparatively small, annotated datasets (i.e., the user-provided
95 videos and labels). This aligns pretrained representations with the specific demands of behavioral
96 segmentation rather than retraining the model from scratch [34] (**Figure 1f**).

97 To generate frame-level predictions, FERAL extends the encoder with an attention-based pooling
98 module and a classification head. Each video chunk is represented as a sequence of spatiotemporal
99 tokens that are integrated by a transformer, then compressed via attention pooling into temporally
100 aligned embeddings. These embeddings are normalized and linearly projected into class logits,
101 producing frame-wise probabilities that are ensembled across overlapping segments.

102 The resulting outputs are interpretable ethograms that describe what actions animals are
103 performing at each frame. These can be visualized directly or integrated into downstream analyses.
104 Model weights are automatically saved, allowing users to reuse or fine-tune trained networks on new
105 datasets without retraining from scratch (**Figure 1g**).

106 Together, these design choices make FERAL both powerful and accessible, establishing a
107 practical and reproducible framework for direct video-to-behavior mapping.

108 **FERAL outperforms state-of-the-art methods across benchmarks**

109 To evaluate FERAL’s precision and robustness relative to existing approaches, we benchmarked
110 its performance on two established datasets that provide raw videos aligned with frame-level
111 behavioral annotations. The first, the Caltech Mouse Social Interactions (CalMS21) dataset [35],
112 contains recordings of freely behaving mice engaged in resident–intruder assays, paired with both
113 tracked poses and expert behavioral labels (**Figure 2a**). The second, MaBE (Multi-Agent Behavior
114 Benchmark) [36], is a large-scale, multi-species dataset spanning mice, beetles, ants, and flies,
115 each annotated across diverse behavioral categories. Together, these datasets test FERAL’s core
116 capability: direct video-to-behavior mapping without reliance on pose trajectories or bounding-box
117 detections. Because only the beetle subset of MaBE provides both raw video and synchronized
118 expert annotations, our evaluation focused on this subset.

119 Originally introduced as a community challenge for the classification of social behavior, the
120 CalMS21 dataset included multiple baseline models and public submissions. We compared FERAL
121 against a suite of alternative pose- and video-based approaches using mean average precision (mAP),
122 the official metric of the original challenge. The strongest released baseline, in addition to the labeled
123 behavioral dataset, employed self-supervised pretraining on large sets of unlabeled pose trajectories
124 [35]. We also report the Competition Top-1 entry from the official leaderboard and results from
125 Google’s VideoPrism paper [37]. This closed-source video-understanding model was fine-tuned on
126 CalMS21 by freezing its backbone and training only the attention-pooling and classification heads.

127 FERAL achieved the highest overall performance on CalMS21, reaching 94.5%, outperforming
128 the strongest baseline (88.9%), competition Top-1 (91.4%), and VideoPrism (91.5%) models
129 (**Figure 2b-c, Supplementary Video 1**). FERAL also maintained consistent performance across
130 videos, correctly classifying the majority of the frames within each sequence (**Figure 2d**). Confusion-
131 matrix analysis (**Figure 2e**) showed that residual errors were mostly confined to confusions
132 between “no label” and “investigate” categories, indicating occasional uncertainty in distinguishing
133 background periods from subtle actions. Quantitative comparison of predicted and annotated total
134 behavior durations per video (**Figure 2f**) confirmed high agreement across all behavioral categories,

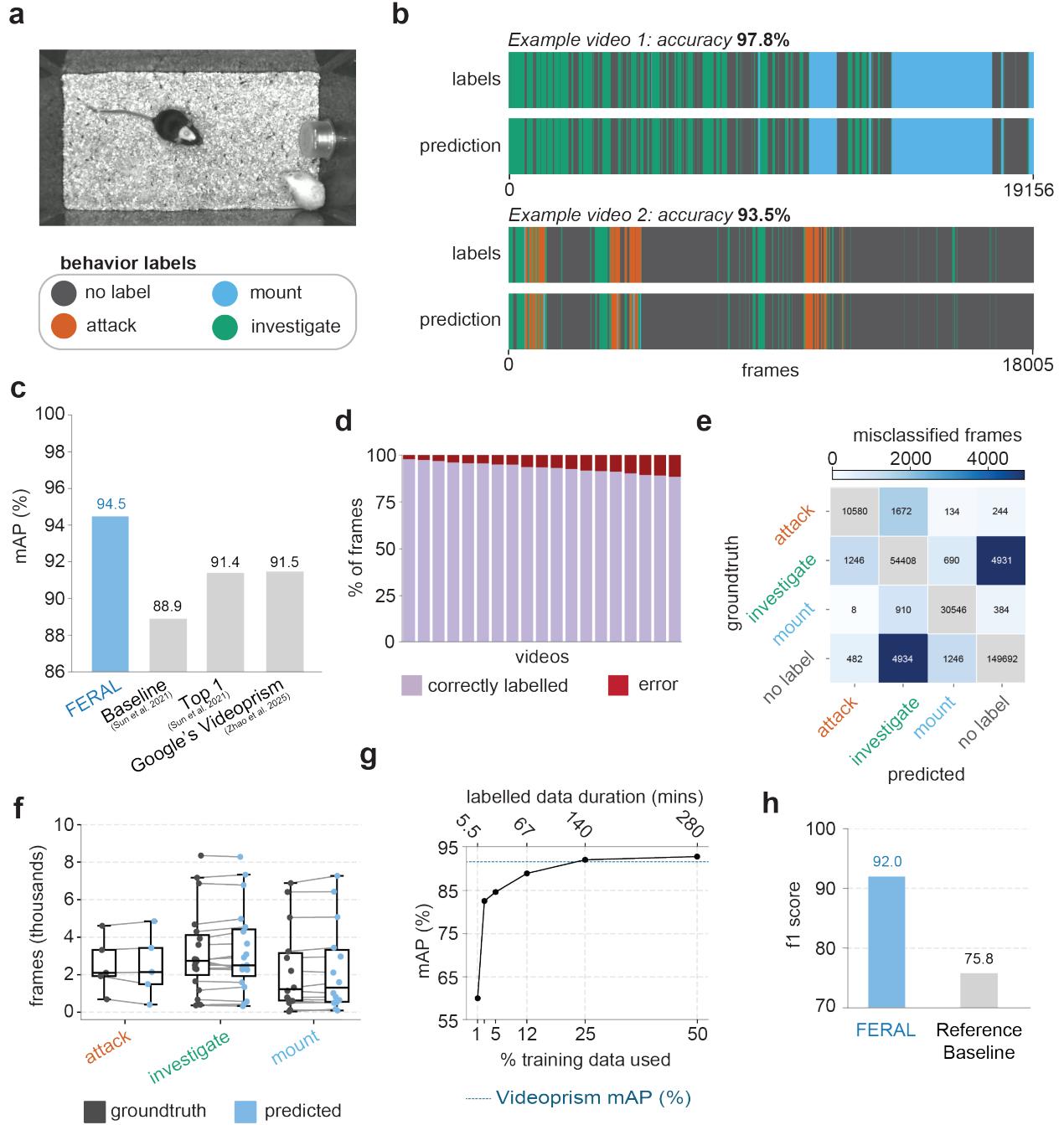


Figure 2: **FERAL outperforms state-of-the-art baselines across benchmarks.** **(a)** Example frame and corresponding behavioral labels from the CalMS21 dataset. **(b)** Representative ethograms comparing FERAL predictions (bottom) to expert annotations (top). **(c)** Mean average precision across models on the CalMS21 dataset. **(d)** Fraction of correctly (lavender) and incorrectly (red) classified frames per video. **(e)** Confusion matrix of misclassified frames. **(f)** Comparison between groundtruth and predicted frame counts for each behavioral category across videos. **(g)** Data-efficiency analysis showing mean average precision (mAP) as a function of the percentage of training data used (1%, 2.5%, 5%, 12%, 25%, 50%). **(h)** F1 score comparison between FERAL and the reference baseline on the MaBE dataset.

135 demonstrating that FERAL accurately preserved both the temporal structure and balance between
136 classes of the annotated behavioral data.

137 To assess how much labeled data FERAL requires to achieve strong performance, we progressively
138 subsampled the training set of CalMS21 (**Figure 2g**). FERAL achieved a mAP of 92.8 using only
139 50% of the available training data and 92.1 with 25%, already surpassing both the VideoPrism and
140 Competition Top-1 models trained on the full dataset. Even when trained with only 12% of the
141 data, performance remained high (89.0%), and robust mAP was maintained down to 5% (84.6%)
142 and 2.5% (82.6%). At the lowest sampling level (1%), FERAL still achieved 60.0%, demonstrating
143 strong data efficiency. These results highlight that FERAL achieves high and stable performance
144 with minimal annotation effort.

145 For the beetle subset of *MaBE*, we evaluated FERAL’s frame-level predictions against expert
146 annotations using the macro-averaged F1 score, following the evaluation protocol of the original
147 publication. The top-performing submission in the *MaBE* challenge achieved a score of 0.758,
148 which we adopt here as the reference baseline [36]. FERAL exceeded this performance by a wide
149 margin (0.92) (**Figure 2h**). Note that in the competition, participants did not have access to labeled
150 data for the target behaviors during training, so it is expected that FERAL, which is trained directly
151 on the target-behavior labels, achieves higher performance.

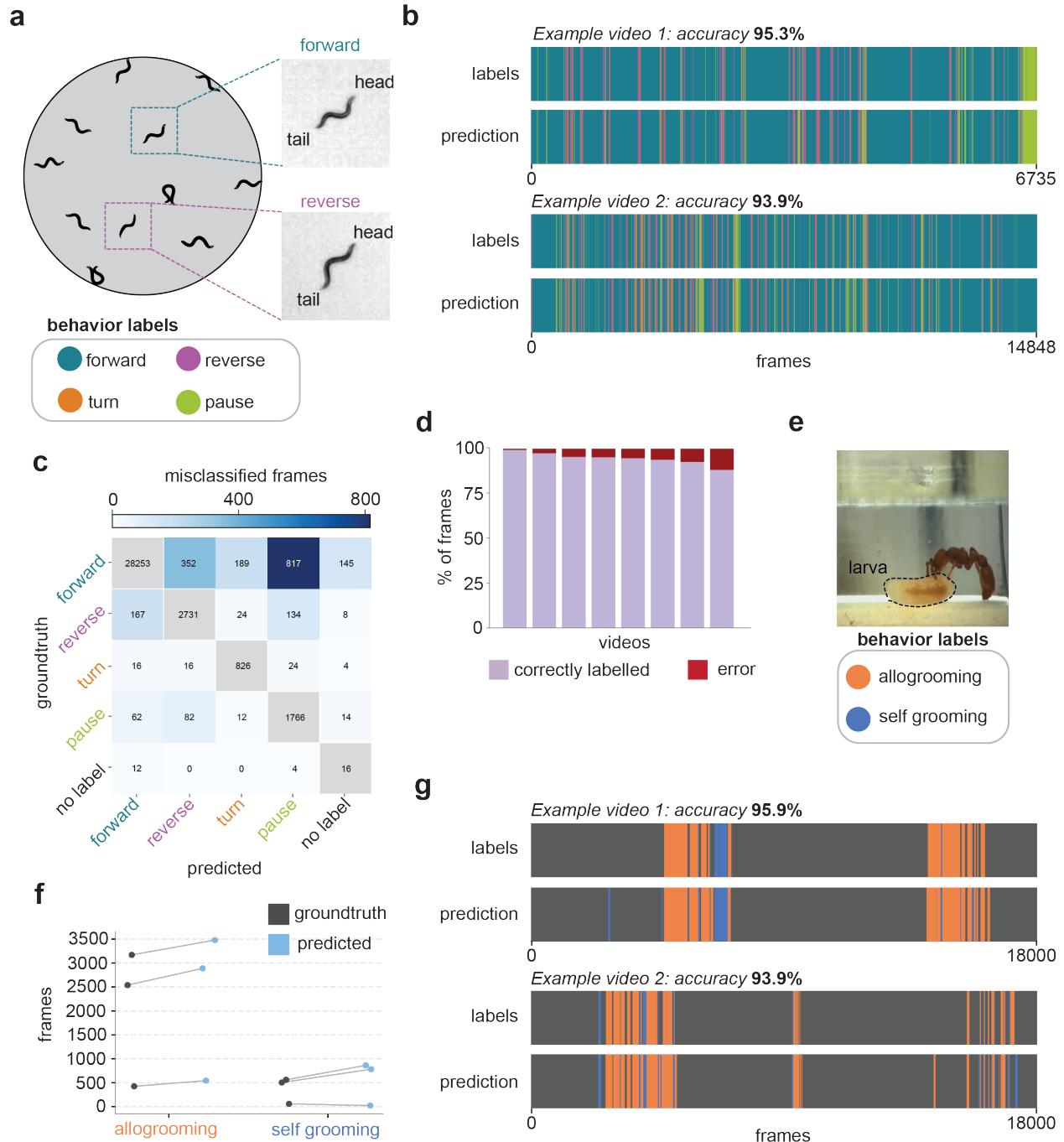


Figure 3: **FERAL captures the temporal and visual appearance of behaviors.** (a) Example frames from the *C. elegans* dataset showing forward and reverse locomotion states. (b) Representative ethograms comparing groundtruth annotations (top) and FERAL predictions (bottom) for *C. elegans* locomotory states. (c) Confusion matrix of misclassified frames across the four locomotor states (*forward, reverse, turn, pause*). (d) Fraction of correctly (lavender) and incorrectly (red) classified frames per video for *C. elegans*. (e) Example frame from recordings of dyadic interactions between an adult and larval clonal raider ant (*Ooceraea biroi*). (f) Quantitative comparison of groundtruth and predicted frame counts for each grooming category. (g) Representative ethograms from *O. biroi* videos comparing groundtruth and FERAL predictions for self-grooming and allogrooming events.

152 **FERAL captures the temporal structure and visual appearance of behaviors**

153 To assess FERAL’s ability to generalize across species and recording modalities, we applied it
154 to two datasets, both acquired in laboratory conditions, representing distinct levels of behavioral
155 organization and scene complexity.

156 **Single-animal behavior in *C. elegans*.** This dataset is comprised of recordings of freely moving
157 *Caenorhabditis elegans* performing four canonical locomotor behaviors: forward crawling, reverse
158 crawling, turning, and pausing (**Figure 3a**) [38, 39]. Unlike other datasets used in this study, these
159 behavioral labels were not manually annotated; instead, each frame was automatically assigned
160 to a locomotory state using a rule-based heuristic pipeline adapted from prior work [38]. In this
161 pipeline, worms were segmented from the background, centroid positions were extracted, head-tail
162 orientation was inferred from midline dynamics, and locomotor state was determined using speed
163 thresholds, direction of motion, and self-intersection criteria (see Methods).

164 For each video, the centroid of the worm was continuously tracked, and a cropped image window
165 centered on the animal was extracted to maintain consistent framing. FERAL leverages temporal
166 context by integrating information from preceding and subsequent frames to classify each moment
167 in time. As a result, it accurately segmented all four behavioral states and correctly distinguished
168 between forward and reverse locomotion. These two locomotory states appear nearly identical in
169 single frames but become separable through their temporal dynamics (**Figure 3b, Supplementary**
170 **Video 2**).

171 Confusion-matrix analysis (**Figure 3c**) showed that the majority of errors were confined to
172 transitions between kinematically adjacent states, such as “forward” versus “pause” or “forward”
173 versus “reverse.” Across videos, FERAL maintained very high frame-level accuracy (**Figure 3d**).
174 Beyond its strong performance, this analysis illustrates how FERAL can be seamlessly integrated into
175 existing behavioral pipelines as a post-tracking module, complementing traditional centroid-based
176 approaches by directly mapping cropped video segments to behaviors.

177 **Dyadic interactions in *Ooceraea biroi*.** We next tested FERAL on a dataset capturing adult–larva
178 interactions in the clonal raider ant *Ooceraea biroi*, which included expert annotations (J.Z.) for
179 self-grooming and allogrooming events (**Figure 3e**). Self-grooming involves individuals cleaning
180 their own body, while allogrooming targets another colony member which, in this case, is a larva.
181 Accurate classification thus requires recognizing not only motion patterns but also the spatial
182 relationship between the adult and the larva. This poses a major challenge for pose-based pipelines:
183 (i) pose estimation is unreliable under frequent occlusions; and (ii) modeling these social interactions
184 using skeletonized data requires integrating information from multiple individuals.

185 FERAL bypasses the need for pose estimation and pose-based segmentation, thus avoiding these
186 challenges altogether. It reliably identified both self- and allogrooming events directly from raw
187 videos (**Supplementary Video 3**), maintaining the overall temporal dynamics of each behavior
188 across videos with high fidelity to the original expert annotations (**Figure 3f-g**).

189 Altogether, these results highlight FERAL’s ability to extract meaningful behaviors from
190 spatiotemporal patterns that are difficult to access with single-frame or pose-based approaches,
191 demonstrating its capacity to generalize from individual locomotion to social behaviors.

192 **FERAL generalizes to field recordings of wild animals**

193 Given FERAL’s robust performance on datasets acquired in the laboratory, we next evaluated its
194 generalization to field recordings, which typically exhibit higher scene complexity than laboratory
195 videos.

196 **Vigilance behavior in zebras.** This dataset consisted of videos of wild Grevy’s (*Equus grevyi*)

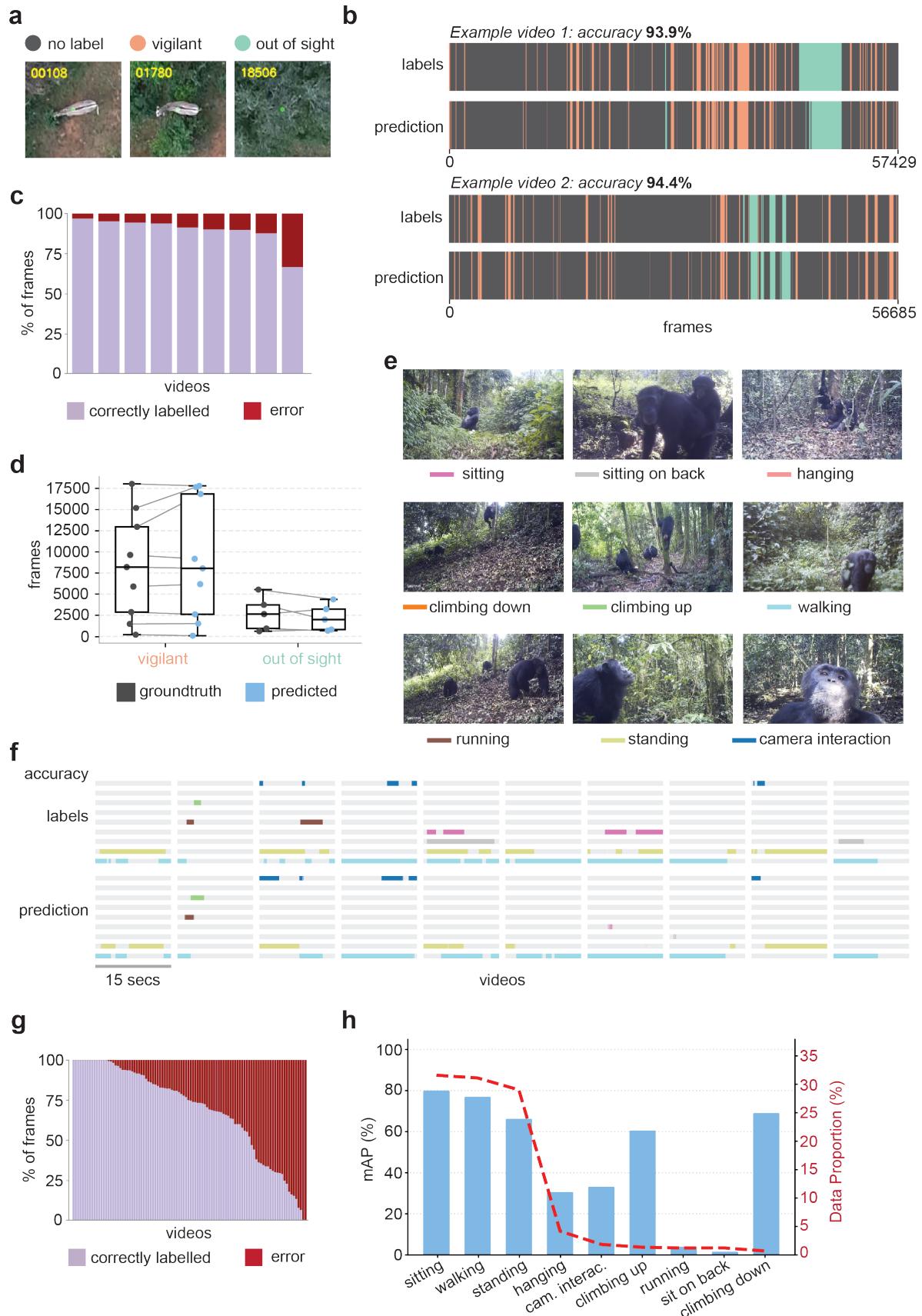


Figure 4: **FERAL generalizes to field recordings of wild animals.** (a) Representative drone frame from recordings of wild Grevy’s zebras (*Equus grevyi*) in Kenya. (b) Example ethograms comparing expert annotations (top) and FERAL predictions (bottom) for the zebras dataset. (c) Fraction of correctly (lavender) and incorrectly (red) classified frames per video for the zebras dataset. (d) Quantitative comparison of predicted and annotated vigilance and out-of-sight bout durations across videos for the zebras dataset. (e) Representative frames from the PanAf500 dataset [41], showing wild chimpanzees recorded by camera traps in forest environments, with one example frame for each annotated behavioral class. (f) Representative ethograms comparing expert labels (top) and FERAL predictions (bottom) across locomotor and postural behaviors for the PanAf500 dataset. (g) Fraction of correctly (lavender) and incorrectly (red) classified frames per video for the PanAf500 dataset. (h) Mean average precision (mAP) per behavioral class (bars) and labeled data proportion (red line) for each behavioral class for the PanAf500 dataset.

197 captured using a drone in Kenya. Free-ranging groups of zebras were filmed from a nadir (directly
198 overhead) perspective. As with the *C. elegans* dataset, the centroid of each animal was continuously
199 tracked following methods described in [40] and a cropped image window was centered on the
200 animal to maintain consistent framing. Individual videos were annotated in BORIS by an expert
201 (B.R.C.) to identify bouts of vigilance behavior, defined as the individual standing still with its
202 head raised (**Figure 4a**). Due to the nadir perspective, the vertical head position of the zebras is
203 challenging to detect, and previous attempts using pose estimation methods were unsuccessful at
204 reliably identifying vigilance bouts.

205 FERAL accurately detected the onset and duration of vigilance periods directly from raw video
206 (**Figure 4b**), closely matching expert annotations across most recordings. FERAL also accurately
207 detected the frames in which the animal was out of sight. The fraction of correctly labeled frames
208 was consistently high across all videos (**Figure 4c**). Quantitative comparison between predicted
209 and annotated frame durations showed strong agreement for both *vigilant* and *out-of-sight* states
210 (**Figure 4d**), indicating that FERAL successfully generalized to aerial perspectives and preserved
211 temporal precision even under the natural variability of field conditions.

212 **Chimpanzee and gorilla behavior in the wild.** To further assess FERAL’s generalization to
213 complex, naturalistic scenes, we evaluated its performance on the *PanAf500* dataset [41], which
214 contains camera-trap videos of wild chimpanzees (*Pan troglodytes*) and gorillas recorded across
215 multiple African field sites as part of the Pan African Programme: The Cultured Chimpanzee.

216 Each clip was manually annotated with fine-grained behavioral categories, including locomotor
217 (e.g., walking, climbing, running), postural (e.g., standing, sitting, hanging), and interaction-related
218 behaviors (**Figure 4e**). Multiple behavior labels can be assigned to the same frame, making this
219 a multi-class, multi-species dataset that tests FERAL’s capacity to handle visual and behavioral
220 complexity.

221 Despite substantial variation in lighting, vegetation, and visibility typical of camera-trap footage,
222 FERAL accurately segmented common postural and locomotor behaviors such as *sitting*, *walking*,
223 *standing*, and *climbing* (both up and down), achieving an overall mean average precision of 65.7%.
224 The fraction of correctly labeled frames, where predictions for all behaviors are correct, varied
225 widely across videos, possibly reflecting differences in scene complexity, visibility, and behavioral
226 diversity (**Figure 4f-g, Supplementary Video 4**). While rare classes such as *running*, and *sitting on*
227 *back* exhibited lower precision, visually distinctive yet infrequent behaviors like *climbing up* and
228 *climbing down* were detected with high average precision despite their low representation in the
229 labeled data (**Figure 4h**).

230 Together, these results demonstrate that FERAL’s performance on videos acquired in laboratory
231 settings transfers robustly to field recordings spanning diverse species, habitats, and acquisition
232 modalities. Its ability to extract meaningful behavioral structure directly from raw video highlights
233 the model’s versatility and scalability for quantifying natural behavior in ecologically realistic
234 contexts.

235 **FERAL identifies emergent collective behaviors directly from raw video**

236 To test whether FERAL generalizes from individual and dyadic interactions to collective dynamics,
237 we applied it to recordings of clonal raider ant (*Ooceraea biroi*) colonies filmed continuously
238 over several days (**Figure 5a**). In this species, foraging occurs as synchronized group raids [42].
239 Manual annotation of such raiding events was used as ground truth. Despite the high density of
240 individuals, FERAL accurately detected the onset and duration of collective raids directly from raw
241 video frames (**Figure 5b**). FERAL captured these events as emergent visual signatures without

242 requiring individual tracking or explicit modeling of group structure. By contrast, conventional
243 multi-animal tracking approaches would need to reconstruct trajectories for each individual [21, 43]
244 and subsequently infer collective states through post-tracking analyses [42]. FERAL bypasses these
245 steps, enabling scalable, tracking-free quantification of collective behavior.

246 Confusion-matrix analysis (**Figure 5c**) revealed that misclassified "raiding" frames were not
247 random but systematically associated with moments of intense movement, when many ants had
248 exited the nest and were exploring the arena. Conversely, most "no raiding" misclassifications
249 occurred during phases of a raid when most ants were in the nest and overall activity declined.
250 These patterns suggest that FERAL's classifications are driven by emergent visual cues, such as
251 colony-level motion and spatial distribution of individuals, demonstrating that the model learns the
252 collective visual signature of raiding behavior directly from scene appearance.

253 By recognizing colony-level behaviors directly from video, FERAL extends the reach of direct
254 video-to-behavior analysis to emergent social dynamics that arise from distributed coordination
255 among many animals. This ability opens new possibilities for studying the neural, ecological, and
256 evolutionary principles governing collective behavior across species and environments.

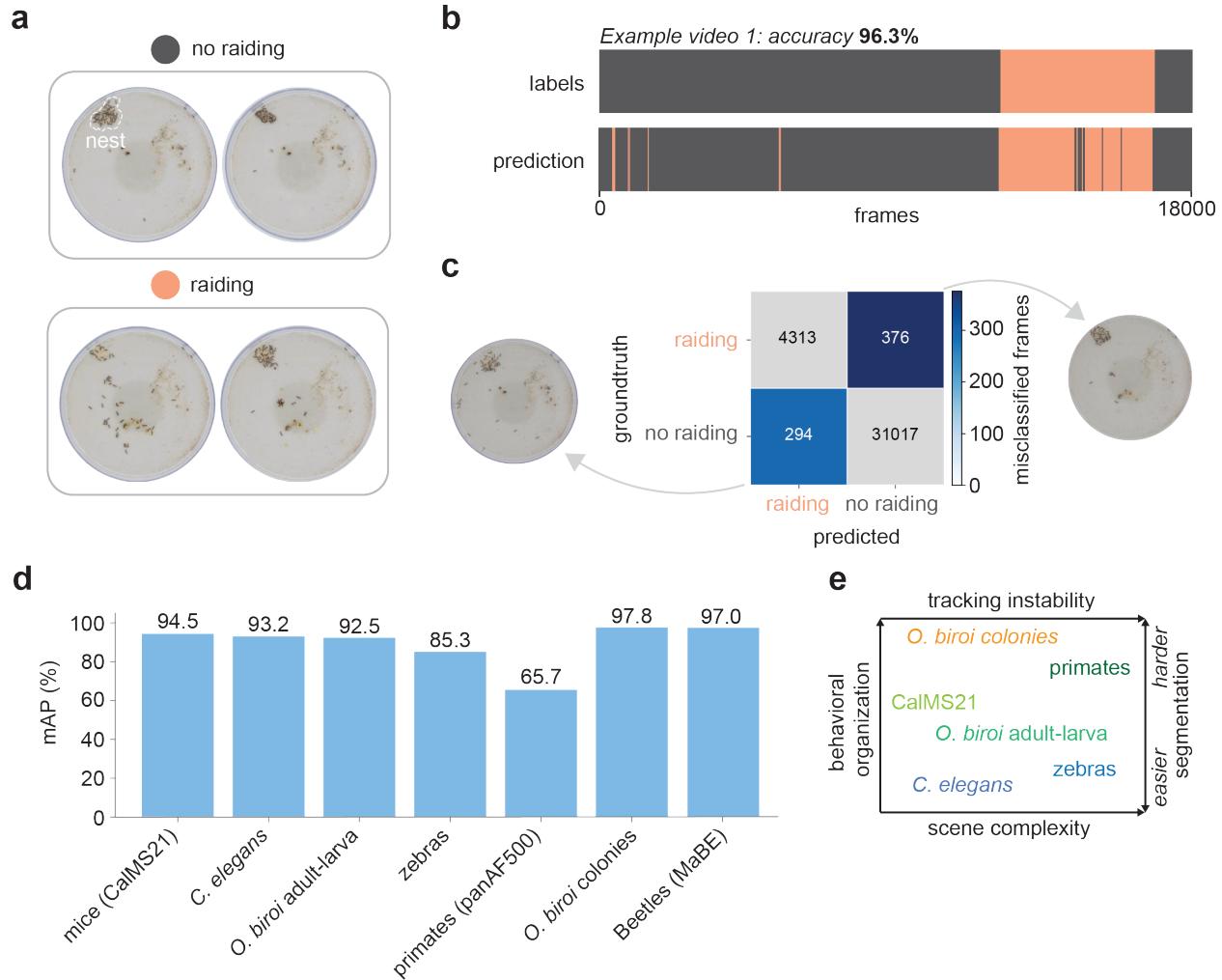


Figure 5: **FERAL captures emergent collective behavior and achieves high performance across dataset.** (a) Representative frames from recordings of clonal raider ant (*Ooceraea biroi*) colonies during non-raiding (top) and raiding (bottom) phases, with nests indicated by dashed outlines. (b) Example ethogram comparing groundtruth annotations (top) and FERAL predictions (bottom). (c) Confusion matrix showing correspondence between groundtruth and predicted frames across all videos. Representative frames showing visual appearance of misclassified frames. (d) mAP of FERAL across datasets tested in this study. (e) Summary of diversity in both behavioral organization and scene complexity across the dataset tested in this study.

257 Discussion

258 The quantitative study of animal behavior has advanced rapidly with the integration of computer
259 vision and deep learning. Existing methods succeed at answering where animals and their body parts
260 are over time, yet determining what actions animals are performing remains challenging, especially
261 for complex environments. This limitation constrains discovery across ethology, neuroscience, and
262 ecology, where understanding behavior and its structure is essential.

263 FERAL addresses this gap by reframing behavioral quantification as a video-understanding
264 problem. Instead of relying on intermediate abstractions such as pose or trajectory, FERAL maps raw
265 video directly to frame-level behavioral labels, capturing both the temporal and visual appearance of
266 animal actions. It achieves high performance across datasets, detecting discrete actions with high
267 precision on a wide range of scene complexities (**Figure 5d-e**). This generalization is enabled by its
268 design. Its foundation-model backbone (e.g., V-JEPA2) allows robust fine-tuning even with limited
269 behavioral annotations [29].

270 FERAL is an open-source, freely available toolkit that includes comprehensive documentation,
271 tutorials, and benchmark datasets (www.getferal.ai), lowering the technical barrier for users
272 across disciplines. Its modular architecture supports all stages of analysis, from preprocessing
273 and training to inference, and can be deployed locally, on high performance clusters, or through
274 cloud-based GPU services such as Google Colab and RunPod.

275 FERAL also integrates easily with existing pipelines. Although it focuses on frame-level

276 segmentation and does not assign persistent identities, it can be paired with multi-animal tracking
277 systems such as TRex [21], which produce animal-centered video segments suitable for direct
278 inference. Together, these tools provide a unified solution to the two central questions of behavioral
279 quantification: where animals are and what actions they are performing.

280 By combining strong performance with user-friendly design, FERAL establishes a robust
281 foundation for behavioral segmentation across species and experimental paradigms. Behavior is
282 central to many disciplines, ranging from neuroscience to ecology, and its quantitative analysis is a
283 common point of entry to understand how biological phenomena work. FERAL gives biologists who
284 study behavior access to the latest video-understanding models, streamlining behavioral analyses
285 and enabling experiments that were previously out of reach, thereby accelerating discovery across
286 fields.

287 **Methods**

288 **Datasets**

289 We evaluated FERAL across seven datasets spanning a range of taxa, behavioral contexts, and
290 recording modalities. All datasets provided raw video data paired with frame-level behavioral
291 annotations, enabling direct video-to-behavior mapping. Other open-source datasets containing only
292 keypoint trajectories or bounding boxes without frame-level behavioral labels were excluded, as the
293 current implementation focuses on frame-level classification. Training and test statistics, including
294 the number of frames per split for each dataset, are reported in Table 2, and class distributions are
295 given in Table 3. Note that for single-class classification, we include an “other” class that the model
296 is trained to predict when no other relevant behaviors are present. For all reported class-averaged
297 metrics, the “other” class is excluded

298 **CalMS21 (Mice social interactions)**

299 The Caltech Mouse Social Interactions (CalMS21) dataset captures freely behaving mice engaged in
300 resident–intruder assays [35]. Each recording includes tracked pose keypoints and frame-level labels
301 for social behaviors such as attack, mount, and investigation. Following prior work, we evaluated
302 performance using the mean average precision (mAP) metric on the same train–test split defined in
303 the original challenge.

304 **MaBE (Multi-Agent Behavior Benchmark)**

305 MaBE [36] is a large-scale multi-species benchmark designed for evaluating multi-agent behavior
306 analysis. It includes recordings of mice, beetles, ants, and flies paired with expert-annotated
307 behavioral categories. As only the beetle subset provides raw video aligned with human frame-level
308 labels, we restricted our evaluation to this partition. Behavioral categories included locomotion,
309 contact, and interaction states, annotated by trained ethologists. We report macro-averaged F1 scores
310 following the evaluation protocol described in the original publication.

311 **PanAf500**

312 PanAf500 [41] is a dataset of wild primate behavior, featuring both chimpanzees and gorillas
313 recorded by camera traps across multiple African field sites within the Pan African Programme.
314 The cameras operate during the day and at night, capturing a range of behaviors expressed by single
315 individuals and groups under diverse weather and lighting conditions.

316 The original work introduces two datasets: PanAf20k, containing 20k clips with a single behavior
317 label per clip, and PanAf500, consisting of 500 clips in which each individual ape is detected with a
318 bounding box in every frame and assigned a behavior label. Because we wanted to evaluate FERAL
319 on per-frame behavior classification, we used the smaller PanAf500 subset. For this dataset, we
320 ignored the bounding box coordinates and used only the behavior labels, pooling behaviors across
321 all bounding boxes in a frame into a single multi-label target (e.g., if two apes were labeled as sitting
322 and one as climbing, the frame-level target would be sitting and climbing).

323 In our experiments, we retained the original training split (400 15s videos) and constructed a
324 single held-out test set by combining the validation and test splits from the original paper (100 15s
325 videos).

326 Because we significantly changed the objective for this dataset, we do not compare our
327 performance to the results reported in the original paper. For PanAf500, the original work reports
328 per-behavior accuracies by predicting a single class for individual behavior sequences for each
329 detected animal, whereas we predict multiple classes for the whole frame at each frame of the
330 video. Beyond the original setup, strong results on this dataset have also been obtained using
331 self-supervised pretraining of a V-JEPA2 model [44]. However that approach follows the original
332 work by first cropping individual animals before applying a separate classifier, while also adding an
333 extra pretraining stage. In contrast, for this dataset we intentionally keep the pipeline simple and
334 evaluate how well a model can perform without any intermediate detection or cropping steps.

335 ***C. elegans* locomotion**

336 Animals were grown at 20 °C on nematode growth media (NGM) plates seeded with *E.coli* OP50
337 bacteria [45]. Experiments were performed on young adult hermaphrodites (picked as L4s 14 -
338 16hrs before).

339 Animals are expressing two transgenes: one extrachromosomal array kyEx = pSM(F23H12.7p:
340 :ReaChR::sl2::GFP) at 70ng/uL + pSM(myo3p::mCherry) at 3ng/uL and one integrated transgene
341 kyIls= = pSM(ser-4p::flp, sto3-p::frt::HisC11::sl2::mCherry, ges-1p::nls-GFP). The animals were not
342 exposed to all-trans-retinol nor histamine during the course of the recordings.

343 Animals were recorded while performing an off-food foraging assay [46, 47], using a plastic
344 ring (6mm Clear Mylar Stencil Sheets) instead of *CuCl*₂ as a boundary to contain the animals.
345 Preconditioning plates were made the night before by seeding NGM plates with a thin uniform OP50
346 lawn 16hrs prior to the start of the assay. 45 minutes prior to the start of the assay, a 1 in x 1 in
347 plastic ring was placed on the preconditioning plate as a boundary and 20 young adult worms were
348 picked onto the lawn. 5 minutes prior to the start of the assay, 8-11 animals were transferred to an

349 unseeded NGM plate to clean off excess food, then transferred to the assay plate, an unseeded 10cm
350 NGM plate with a plastic ring (2 in x 2 in) boundary to keep animals within the recording field of
351 view. Behavior was recorded for 45 min using a Basler ace acA5472-17 μm USB 3.0 Monochrome
352 Camera at 6 fps and 3,648px x 3,648px FOV.

353 Animals were tracked using custom python software which segments the worms from background
354 using thresholding and generates tracks of the centroid of segmented worms using intersection over
355 union. This results in tracks including centroid and mask of the worm.

356 Behavior was classified heuristically using rules adapted from [48, 49] implemented in custom
357 python software. For each frame, it was determined if the animal was self-intersecting by using
358 the ratio of the area: perimeter ratio of the outermost contour of the mask (ratio>3.3). For non
359 self-intersecting frames, masks were skeletonized to get the midline of the worm and aligned in the
360 same direction by minimizing the distance between midline points in adjacent frames. The head-tail
361 vector was taken from the first and last midpoints of the worm.

362 Velocity was calculated by taking the difference between centroids 4 frames (=0.67s) apart and
363 dividing by $dt = 0.67\text{s}$. Head versus tail was assigned as the direction along the head-tail axis the
364 animal moves in more often (the animal moves towards the direction of its head more often than the
365 direction of its tail). Speed was calculated by taking the absolute magnitude of the velocity. Speed
366 was signed by the direction of the velocity relative to the direction of the head-tail vector.

367 Turning was classified as frames in which either a) both i) midpoint-tail vector (the vector
368 connecting the midpoint and the tail of the worm) was greater in length than the midpoint-head
369 vector (the vector connecting the midpoint to the head of the worm) and ii) the angle between the
370 midpoint-tail vector and midpoint-head vector was less than 45 degrees or b) when the animal was
371 self-intersecting.

372 Pausing was classified as frames for which speed < 50 $\mu\text{m}/\text{s}$ for at least 0.5s and for which the
373 animal was not turning.

374 Forwards was classified as speed greater than 0 $\mu\text{m}/\text{s}$ and not turning or pausing.

375 Reversal was classified as speed less than 0 $\mu\text{m}/\text{s}$ and not turning or pausing.

376 **Collective behavior and adult-larva interactions in clonal raider ant (*Ooceraea biroi*)**

377 Stock colonies of *Ooceraea biroi* were maintained in constant light at 25 °C in Tupperware containers
378 (40 x 26 cm) with a 2 cm thick plaster of Paris floor. Colonies were fed with frozen fire ant (*Solenopsis*
379 *invicta*) brood following the lab's regular feeding schedule (3 times per week) and cleaned and
380 watered once per week, as needed.

381 For behavior experiments, adult ants and fourth-instar larvae from the same stock colony (clonal
382 line B genetic background) were collected and housed in 5 cm Petri dishes lined with a plaster of
383 Paris floor and kept at 25 °C. Behavioral assays were performed in a custom-built acrylic chamber
384 with transparent sides and plaster of Paris floor. In each assay, an adult ant was allowed to settle in
385 the chamber for 30 minutes before a larva was introduced. Videos were recorded from the side at
386 10 frames per second through a FLIR blackfly camera (BFS-U3-50S5C-C) and lens (Computar,
387 MLM3X-MP), using Spinnaker Software Development Kit, at a pixel resolution of 2448px x 2048px.
388 Recordings were collected across three consecutive days.

389 Behavioral annotations of adult ants and larvae were performed manually using ELAN [50]. We
390 scored adult self-grooming and adult allogrooming of larvae, and recorded the start and end times
391 and duration of these behaviors (in milliseconds).

392 For the colony tracking experiment, groups of 100 adults and 100 larvae were randomly
393 subsampled from stock colonies and established in Petri dishes (90 x 20mm) lined with a humidified
394 plaster of Paris base. Sub-colonies were allowed to acclimate for four days. Sub-colonies were
395 then continually video recorded at 0.1 fps for 30 days under constant illumination at 25 °C.
396 Approximately every two days colonies were fed with frozen *Solenopsis invicta* brood and were
397 cleaned approximately every two days.

398 **Zebras recordings from drones**

399 Herds of Grevy's zebra were filmed at Mpala Conservancy in Laikipia, Kenya in 2017 and 2018
400 using DJI Phantom 4 Pro Drones (DJI, Shenzhen, China). Drone import and operations were
401 authorized by the Kenya Civil Aviation Authority (KCAA) and carried out by a licensed pilot assisted

402 by an observer who maintained visual contact with the drone and a ground observer who maintained
403 situational awareness. During filming, the drones were positioned directly above the group at a
404 height of approximately 85 m above ground level (AGL), and followed the group's movements. To
405 achieve continuous observations longer than the drone's battery duration, two drones were flown
406 in a relay: when one drone's battery became depleted, a second drone was positioned 10 m above
407 the first one. The first drone was then recalled to the launch point and the second drone lowered
408 down 10 m to continue following the group. If animals seemed disturbed by the drone (e.g. running
409 away) or moved too far from the launch point, the observation was terminated. When groups were
410 calm and remained in range, observations spanned 3 drone flights (approximately 45-50 minutes).
411 During the first two flights, groups were filmed in an undisturbed state. During the third flight,
412 researchers approached the group on foot to elicit a detection and flight response, during which the
413 drone followed the group until they ran out of range. All fieldwork in Kenya was conducted with the
414 permission of the National Commission for Science, Technology and Innovation and in affiliation
415 with the Kenya Wildlife Service. Data collection protocols were reviewed and approved by Ethikrat,
416 the independent ethics council of the Max Planck Society.

417 Drone recordings were initially captured at 4K resolution and 60 fps, but were downsampled
418 to 30 fps prior to further processing. A multi-stage pipeline was applied to generate continuous
419 movement trajectories for all animals in the recordings [40]. Recordings were then cropped to
420 generate individual videos of a small square area (160 x 160 px – 210 x 210 px) centered on each
421 animal. Forty-five individual videos from four observations were then manually annotated in BORIS
422 [31] to identify bouts of vigilance behavior, defined as the animal standing still with its head raised,
423 and periods when the animal was out of sight, for example due to passing under occluding vegetation.
424 All other video frames were uncategorized.

425 Model Architecture

426 FERAL fine-tunes a state-of-the-art video-understanding backbone to perform frame-level behavioral
427 classification across predefined categories. The model outputs class probabilities for every frame.

428 The complete default configuration used for all reported dataset, including all hyperparameters, is
429 available on GitHub¹

430 **Backbone evaluation and selection**

431 We systematically evaluated several recent video-understanding architectures as potential encoders
432 for FERAL, focusing on their balance of accuracy, computational efficiency, and ease of deployment.

433 **InternVideo2.** InternVideo2 is a 1-billion-parameter transformer model trained in multiple
434 pretraining stages [51]. Despite strong representational capacity, fine-tuning proved prohibitively
435 expensive, requiring up to eight hours for a full run on the CalMS21 dataset using four H100 GPUs.
436 Furthermore, the codebase depended on numerous bespoke modules, complicating installation and
437 development. While performance was promising, these practical limitations rendered InternVideo2
438 unsuitable for a user-friendly behavioral analysis framework.

439 **SmallVLM2.** We next assessed SmallVLM2, a multimodal video–language model available in
440 configurations ranging from 256M to 2.2B parameters [52]. Although integration via Hugging Face
441 Transformers greatly simplified deployment, performance lagged behind state-of-the-art methods.
442 The model processed frames largely in isolation, aggregating temporal information only in a
443 shallow pooling layer. This architectural constraint, compounded by the predominance of non-video
444 modalities during pretraining, limited its ability to capture long-range motion dynamics. Even
445 with extensive regularization (including partial freezing, data augmentations, and label smoothing)
446 SmallVLM2 exhibited overfitting on small datasets and failed to generalize on internal benchmarks.

447 **V-JEPA2.** Based on these observations, we adopted V-JEPA2 as the encoder [29]. Unlike
448 video–language models, V-JEPA2 is a dedicated video foundation model trained self-supervised
449 on over one million hours of unlabeled video using a masked prediction objective. We employed
450 the smallest available configuration (330M parameters) finetuned on the Diving48 dataset[53],
451 which offers a favorable trade-off between performance and computational cost. Fine-tuning
452 enables FERAL to align pretrained spatiotemporal representations with the specific requirements of

¹https://github.com/Skovorp/feral/blob/main/configs/default_vjepa.yaml

453 behavioral segmentation, achieving great performance across benchmarks with modest data volumes
454 and standard GPU resources.

455 **Classification head**

456 To convert spatiotemporal embeddings from the backbone into frame-level behavioral predictions,
457 we designed a lightweight classification head that aggregates contextual information and outputs
458 per-frame logits across behavioral classes.

459 Each input video chunk is represented as a sequence of thousands of spatiotemporal tokens.
460 These tokens are first processed through transformer layers of the V-JEPA2 encoder[54] that enrich
461 each token with contextual information from the entire sequence. To map this long sequence onto the
462 temporal resolution of the input video, we employ an attention-based pooling module. Specifically,
463 a set of 64 learnable query embeddings cross-attend to the encoder outputs, extracting features that
464 correspond to individual frames.

465 The resulting pooled embeddings are flattened and passed through a Batch Normalization layer
466 [55], which stabilizes training and controls feature variance, followed by a dropout layer ($p = 0.5$)
467 [56] to reduce overfitting. A final linear projection maps the normalized embeddings to class logits
468 for each frame.

469 **Loss function**

470 FERAL employs different loss functions for single-label and multi-label classification. For single
471 label, FERAL is trained using a cross-entropy loss computed at the frame level. To improve
472 generalization and mitigate overconfidence, we apply label smoothing with a factor of 0.1, which
473 encourages the model to distribute probability mass across semantically related classes rather than
474 assigning absolute certainty to a single label.

475 Because behavioral datasets often exhibit pronounced class imbalance, particularly between
476 dominant "background" states and rare but biologically meaningful actions, we incorporate class-
477 specific weighting into the loss. We found that scaling weights by the square root of the inverse class

478 frequency yielded better performance amplifying the contribution of underrepresented behaviors
479 without overcompensating, compared to using inverse-frequency weighting.

480 In the multi-label setting, FERAL is trained with binary cross-entropy loss, weighting each class
481 by $\sqrt{\frac{N_{\text{positives}}}{N_{\text{negatives}}}}$ to emphasize rare classes.

482 We additionally evaluated focal loss, which dynamically down-weights easy examples to focus
483 learning on difficult cases, but found no consistent improvement across benchmarks.

484 **Chunking strategy**

485 Transformer-based video encoders compute pairwise attention across all spatiotemporal tokens,
486 causing computational complexity to scale quadratically with the number of input tokens. As a
487 result, processing entire behavioral recordings end-to-end is infeasible. Consequently FERAL
488 divides each video into overlapping segments, or *chunks*, of fixed length before processing. Each
489 chunk comprises 64 consecutive frames, resized to 256×256 pixels.

490 To capture fine-grained behavioral dynamics, consecutive chunks overlap by 50% (i.e., stride
491 = 32 frames). This design ensures that short behaviors spanning chunk boundaries remain fully
492 visible within at least one receptive field. During training, overlapping windows also increase the
493 effective number of training samples further improving quality.

494 During inference, as the model outputs per-frame predictions for each chunk, we employ a
495 unified post-processing pipeline to gather the final frame-level predictions. If a frame received
496 multiple predictions, we averaged the corresponding class probabilities. For frames without a direct
497 prediction, we linearly interpolated probabilities from the two nearest frames with predictions. This
498 approach effectively ensembles model’s own local predictions.

499 We benchmarked multiple configurations varying both stride and sampling rate (**Extended**
500 **Figure 1**) and found that sampling every frame with 50% overlap yielded the best balance between
501 accuracy and training speed.

502 **Augmentations**

503 To improve generalization across diverse lighting conditions, species, and recording setups, FERAL
504 employs a combination of video and label-space augmentations. For visual augmentation, we
505 adopted *TrivialAugment* [57], which randomly samples from a set of standard image transformations
506 (e.g., brightness, contrast, rotation, color jitter) and applies them at varying strengths. The same
507 augmentation was applied consistently across all frames within a video, preserving temporal
508 coherence while introducing diversity across video samples.

509 In addition, we applied *MixUp* regularization at the batch level [58]. Each augmented sample was
510 formed as a convex combination of two videos, X and \tilde{X} , and their corresponding label sequences, y
511 and \tilde{y} :

512
$$X_{\text{new}} = \alpha X + (1 - \alpha)\tilde{X}, \quad y_{\text{new}} = \alpha y + (1 - \alpha)\tilde{y},$$

513 where $\alpha \sim \text{Beta}(\lambda, \lambda)$. Because FERAL operates at frame resolution, label mixing was
514 performed element-wise across the temporal dimension. This strategy helps to mitigate overfitting
515 on small datasets.

516 **Training**

517 We use the Adam optimizer [59] with a relatively strong weight decay (0.1) to counteract overfitting
518 given the modest size of behavioral datasets. The learning rate follows a schedule consisting of a
519 linear warm-up for the first 20% of iterations, followed by cosine decay. Models are trained for 10
520 epochs, which we found sufficient for convergence across benchmarks.

521 Unlike approaches such as VideoPrism, which freeze the backbone and train only shallow
522 classifiers, FERAL fine-tunes the last 12 out of 24 transformer layers in VJEPAP2, aligning high-level
523 spatiotemporal embeddings with behavioral structure.

524 To further support generalization and out-of-distribution performance, we allow users to enable
525 an optional weight-averaging feature. When activated, FERAL will maintain an exponential moving

526 average (EMA) θ_{EMA} of the model weights θ , which at step t is computed as [60]:

527
$$\theta_{t+1}^{\text{EMA}} = \beta\theta_t^{\text{EMA}} + (1 - \beta)\theta_t,$$

528 with $\beta = 0.999$. During evaluation, both the standard and EMA-weighted models are assessed
529 and users can select the best performing option for their application. In our reported experiments,
530 we do not report metrics from EMA checkpoints, as some of the datasets are relatively small and we
531 preferred to keep the evaluation protocol simple rather than risk over-optimizing on them.

532 Experiments

533 Data efficiency

534 We tested FERAL’s data efficiency by training on smaller CalMS21 subsets and measuring
535 performance drop versus the full-data baseline, while keeping the test set and all hyperparameters
536 the same. We implemented two complementary subsampling schemes:

537 **(1) Video-level subsampling.** We randomly sampled subsets of training videos at 50% (mAP
538 92.0%) and 25% (mAP 93.0%) and trained FERAL on these reduced sets. Because individual
539 recordings vary in length and behavioral composition, smaller subsets (<25%) produced high
540 variance across runs: some samples didn’t have all the classes present and the total number of frames
541 varied substantially. We therefore do not report results below 25%.

542 **(2) Chunk-level subsampling.** To probe sample efficiency under more balanced class dis-
543 tributions, we first processed the full training set into ready-to-train chunks and then randomly
544 subsampled from them. This design mitigates class imbalance and enabled evaluation on much
545 smaller training sets, down to 1% of the original data.

546 Since expanding datasets typically involves annotating additional videos, video-level subsampling
547 best reflects realistic scaling. However, only chunk-level subsampling allows evaluation under
548 extreme reductions without manual video selection. Under both regimes, training on 25% of the data
549 still exceeded the prior SOTA, and other reduced-data settings maintained strong performance. These

550 results indicate that FERAL’s foundation-model backbone and fine-tuning strategy confer substantial
551 sample efficiency, enabling high performance with limited labeled footage. This is especially
552 valuable in behavioral research, where manual annotation is expensive and time-consuming

553 **Chunking strategies**

554 We systematically benchmarked chunking strategies to balance computational efficiency and quality.

555 Each configuration is defined by two parameters:

556 • **Frame stride**: the interval between frames within a chunk (e.g., every frame, every second
557 frame). Larger strides expand the effective temporal window but reduce temporal resolution.

558 • **Chunk overlap**: the proportion of frames shared between consecutive chunks (e.g., 0%, 50%,
559 75%). At 0% each frame in the video appears in only one chunk, at 50% in two, at 66% in
560 three, etc. Greater overlap increases computation but improves contextual continuity and
561 yields multiple predictions per frame that can be ensembled to enhance quality.

562 Performance improved monotonically as stride decreased, with the best results at stride 1
563 (sampling every frame; labeled as "dense" on **Extended Figure 1a**). Adding overlap yielded
564 substantial quality gains, with 66% overlap producing the highest mean average precision (95.0%).
565 Notably, even sparse settings (0% overlap, sampling every fourth frame) exceeded the competition
566 baseline while requiring 8 times less steps than our default configuration.

567 To test whether improvements were simply due to more training steps, we matched the total
568 number of optimization steps by training the base configuration (50% overlap, dense sampling; blue
569 bars on **Extended Figure 1a**) for fewer or more epochs. Dense configuration were clearly superior
570 to sampling every other frame, and moderate overlap was beneficial, although returns diminished at
571 higher overlaps. We therefore adopted 50% overlap with full-frame sampling as the default.

572 **Freezing strategies**

573 To evaluate the impact of layer freezing on performance, we incrementally froze six-layer blocks of
574 the 24-layer V-JEPA2 encoder, proceeding from the input forward. Partial freezing (up to 12 layers)
575 slightly improved test quality, consistent with mild regularization. In contrast, freezing the entire
576 encoder markedly reduced performance (mAP of 88.4%). These results indicate that fine-tuning at
577 least the final layers is necessary to align pretrained features with behavioral segmentation tasks
578 (**Extended Figure 1b**).

579 **Metrics**

580 We report mean average precision (mAP) and macro-averaged F1, precision, and recall. mAP
581 is computed from calibrated probabilities, excluding the “other” category from averaging. For
582 multi-label tasks (e.g., MaBE), results are reported using a fixed threshold of 0.85 for all classes.
583 You can find the performance metrics of FERAL across all benchmark datasets using the same
584 default configuration in Table 1.

585 **Metric rationale**

586 Accuracy was not used, as behavioral datasets typically exhibit substantial class imbalance. Datasets
587 often include a frequent “other” class, while biologically interesting events are relatively rare (often
588 < 1% of frames). In such settings, accuracy overestimates performance by rewarding majority-
589 class predictions. Instead, we report mAP, which offers a more sensitive measure of per-class
590 discrimination across thresholds, while macro F1, precision, and recall provide a complementary
591 view of performance. mAP is computed from calibrated probabilities over the full dataset, whereas
592 the F1 score is computed after thresholding continuous model outputs into discrete predictions.
593 For all reported averages, the “other” class was excluded, as it denotes the absence of annotated
594 behaviors rather than a discrete category.

595 Model selection

596 For all runs, we train the model on the full training dataset and evaluate on the test dataset afterward.
597 We observed that splitting the available non-test data into training and validation sets and selecting
598 the best-performing checkpoint across epochs to evaluate on the test set led to highly variable results,
599 due to the small size of some datasets. Therefore, we chose this setup, accepting that mild overfitting
600 on some datasets is preferable to undertraining on others.

601 Run tracking

602 Each training run automatically logs all metrics to Weights & Biases (W&B), providing users with
603 real-time visualization and reproducibility. Logged outputs include low-level training diagnostics
604 (e.g., loss, learning rate), per-class average precision at both chunk and frame levels, aggregated
605 mAP, and qualitative ethograms after each validation epoch. Each run generates a unique dashboard
606 URL, enabling easy sharing and comparison.

607 Engineering and deployment

608 **Video preprocessing and seekability.** Training requires fast access to arbitrary frames in the
609 video. Many common codecs support only sequential reads and do not allow frame-level seeking,
610 so re-encoding is often necessary. Pre-resizing can also help with HD videos, as decoding large
611 videos may be too slow compared with GPU computation speed. FERAL includes a cross-platform
612 utility that converts input videos to seekable formats and resizes them to 256×256. If source videos
613 are not high resolution and are already seekable, users may skip re-encoding; resizing will then be
614 performed dynamically during training

615 **Hugging Face integration.** To ensure modularity and ease of extension, FERAL leverages
616 the Hugging Face `transformers` library [61] for model management. This abstraction simplifies
617 switching between backbones and streamlines installation compared to earlier architectures (e.g.,
618 `InternVideo2`), which relied on bespoke code and complex dependencies.

619 **Training visualization and monitoring.** FERAL integrates with Weights & Biases (W&B)
620 for experiment tracking, allowing users to monitor training, validation and test performance. Logs
621 are also stored in the cloud, enabling training across multiple servers while maintaining a unified
622 analysis interface. Users can connect their own W&B accounts or use a public workspace provided
623 by FERAL.

624 **Deployment and compute requirements.** FERAL requires 24 GB of VRAM for training,
625 which fits on high-end consumer GPUs. All reported experiments were conducted on NVIDIA H100
626 and L40s GPUs in a university cluster. For users without access to high-performance computing,
627 FERAL provides step-by-step deployment guides for GPU cloud platforms (e.g., RunPod) and a
628 Google Colab notebook, enabling full training runs at low cost without complex setup. All necessary
629 dependencies, including CUDA and PyTorch [62], are pre-installed on these instances, allowing
630 users to begin training within minutes.

631 **Author Contributions**

632 P.S. and J.R. conceptualized the study, developed the method, designed experiments, performed
633 analyses, and co-wrote the manuscript. B.R.C., B.K. and I.D.C. provided the zebra dataset and
634 B.R.C. supported the application of FERAL to the zebra dataset. F.B. provided the *C. elegans*
635 dataset. J.Z. and D.D.F provided the ant grooming dataset. T.K. provided the ant colony dataset.
636 All authors provided feedback on the manuscript.

637 **Data and Code Availability**

638 All source code, training configurations, and documentation for FERAL are publicly available
639 at www.getferal.ai and <https://github.com/Skovorp/feral>. The website and the repository include
640 example datasets, preprocessing utilities, and instructions for local and cloud-based deployment.
641 Data and Supplementary Videos 1-4 used in this publication are available at the GitHub Data
642 Repository (https://github.com/Skovorp/feral_share_data/tree/main)

643 Competing Interests

644 The authors declare no competing interests.

645 Acknowledgments

646 We thank Leslie B. Vosshall and the members of the Vosshall Lab and Data Science Platform at
647 The Rockefeller University, as well as Orli Snir, Giulio Formenti, Yohann Chemtob and Emily
648 B. L. Wright for discussion and comments on the manuscript; This work was supported by the
649 Howard Hughes Medical Institute, and graduate fellowships from the Price Family Center for
650 the Social Brain and the Boehringer Ingelheim Fonds (J.R.). This work was supported by the
651 Rockefeller University, and the Jonathan and Maya Nelson Center for Artificial Intelligence (P.S.).
652 We gratefully acknowledge the Data Science Platform (DSP) at The Rockefeller University for
653 access to the DSP Cluster and computing equipment used in this study. This work was supported
654 by the National Institute on Deafness and Other Communication Disorders under award number
655 K99DC021506 to D.D.F. B.R.C. received support from the European Union’s Horizon 2020 research
656 and innovation program under the Marie Skłodowska-Curie grant agreement No. 748549. B.R.C.
657 acknowledges support from the University of Konstanz’s Investment Grant program. B.K., I.D.C, and
658 B.R.C. acknowledge support from the Deutsche Forschungsgemeinschaft (DFG, German Research
659 Foundation) under Germany’s Excellence Strategy—‘Centre for the Advanced Study of Collective
660 Behaviour’ EXC 2117-422037984, DFG project number 462886202, the European Union’s Horizon
661 2020 Research and Innovation Programme under Marie Skłodowska-Curie Grant 860949, the DFG
662 Gottfried Wilhelm Leibniz Prize 2022 584/22, and the PathFinder European Innovation Council
663 Work Programme 101098722. B.R.C. acknowledges support from NVIDIA Corporation’s Academic
664 Hardware Grant Program. F.B. was supported by a Medical Scientist Training Program grant
665 from the National Institute of General Medical Sciences of the National Institutes of Health under
666 award number T32GM152349 to the Weill Cornell/Rockefeller/Sloan Kettering Tri-Institutional
667 MD-PhD Program and by an F31 Predoctoral Fellowship from the National Institute of Neurological

668 Disorders and Stroke of the National Institutes of Health under award number F31NS132477.

669 References

- 670 [1] André E. X. Brown and Benjamin de Bivort. “Ethology as a physical science”. In: *Nature
671 Physics* 14.7 (2018), pp. 653–657. issn: 1745-2473. doi: 10.1038/s41567-018-0093-0.
- 672 [2] Gordon J. Berman. “Measuring behavior across scales”. In: *BMC Biology* 16.1 (2018), p. 23.
673 doi: 10.1186/s12915-018-0494-7.
- 674 [3] Sandeep Robert Datta et al. “Computational Neuroethology: A Call to Action”. In: *Neuron*
675 104.1 (2019), pp. 11–24. issn: 0896-6273. doi: 10.1016/j.neuron.2019.09.038.
- 676 [4] S E Roian Egnor and Kristin Branson. “Computational Analysis of Behavior.” In: *Annual
677 review of neuroscience* 39.1 (2016), pp. 217–36. issn: 0147-006X. doi: 10.1146/annurev-
678 neuro-070815-013845.
- 679 [5] Marian Stamp Dawkins, Paul Martin, and Patrick Bateson. “Measuring Behaviour. An
680 Introductory Guide”. In: *The Journal of Animal Ecology* 63.3 (1994), p. 746. issn: 0021-8790.
681 doi: 10.2307/5248.
- 682 [6] Konrad Lorenz. “Der Kumpan in der Umwelt des Vogels”. In: *Journal für Ornithologie* 83.3
683 (1935), pp. 289–413. issn: 0021-8375. doi: 10.1007/bf01905572.
- 684 [7] Niko Tinbergen. *The Study of Instinct*. Oxford University Press, 1951.
- 685 [8] Karl von Frisch. “Decoding the Language of the Bee”. In: *Science* 185.4152 (1974), pp. 663–
686 668. issn: 0036-8075. doi: 10.1126/science.185.4152.663.
- 687 [9] Jeanne Altmann. “Observational Study of Behavior: Sampling Methods”. In: *Behaviour*
688 49.3-4 (1974), pp. 227–266. issn: 0005-7959. doi: 10.1163/156853974x00534.
- 689 [10] David J. Anderson and Pietro Perona. “Toward a Science of Computational Ethology”. In:
690 *Neuron* 84.1 (2014), pp. 18–31. issn: 0896-6273. doi: 10.1016/j.neuron.2014.09.005.

- 691 [11] John W. Krakauer et al. “Neuroscience Needs Behavior: Correcting a Reductionist Bias”. In:
692 *Neuron* 93.3 (2017), pp. 480–490. ISSN: 0896-6273. doi: 10.1016/j.neuron.2016.12.041.
- 693 [12] Alex Gomez-Marin et al. “Big behavioral data: psychology, ethology and the foundations of
694 neuroscience”. In: *Nature Neuroscience* 17.11 (2014), pp. 1455–1462. ISSN: 1097-6256. doi:
695 10.1038/nn.3812.
- 696 [13] Jane Goodall. “Tool-Using and Aimed Throwing in a Community of Free-Living Chim-
697 panzees”. In: *Nature* 201.4926 (1964), pp. 1264–1266. ISSN: 0028-0836. doi: 10.1038/
698 2011264a0.
- 699 [14] Margaret Bastock and Aubrey Manning. “The Courtship of *Drosophila Melanogaster*”. In:
700 *Behaviour* 8.1 (1955), pp. 85–110. ISSN: 0005-7959. doi: 10.1163/156853955x00184.
- 701 [15] Cristina Segalin et al. “The Mouse Action Recognition System (MARS) software pipeline
702 for automated analysis of social behaviors in mice”. In: *eLife* 10 (2021), e63720. doi:
703 10.7554/elife.63720.
- 704 [16] Devis Tuia et al. “Perspectives in machine learning for wildlife conservation”. In: *Nature*
705 *Communications* 13.1 (2022), p. 792. doi: 10.1038/s41467-022-27980-y. eprint: 2110.12951.
- 706 [17] Megan Tjandrasuwita et al. “Interpreting Expert Annotation Differences in Animal Behavior”.
707 In: *arXiv* (2021). doi: 10.48550/arxiv.2106.06114. eprint: 2106.06114.
- 708 [18] Alexander Mathis et al. “DeepLabCut: markerless pose estimation of user-defined body parts
709 with deep learning”. In: *Nature Neuroscience* 21.9 (2018), pp. 1281–1289. ISSN: 1097-6256.
710 doi: 10.1038/s41593-018-0209-y.
- 711 [19] Talmo D Pereira et al. “SLEAP: A deep learning system for multi-animal pose tracking.”
712 In: *Nature methods* 19.4 (2021), pp. 486–495. ISSN: 1548-7091. doi: 10.1038/s41592-022-
713 01426-1.
- 714 [20] Dan Biderman et al. “Lightning Pose: improved animal pose estimation via semi-supervised
715 learning, Bayesian ensembling and cloud-native open-source tools”. In: *Nature Methods* 21.7
716 (2024), pp. 1316–1328. ISSN: 1548-7091. doi: 10.1038/s41592-024-02319-1.

- 717 [21] Tristan Walter and Iain D Couzin. “TRex, a fast multi-animal tracking system with markerless
718 identification, and 2D estimation of posture and visual fields.” In: *eLife* 10 (2020), e64000.
719 doi: 10.7554/elife.64000.
- 720 [22] Jacob M Graving et al. “DeepPoseKit, a software toolkit for fast and robust animal pose
721 estimation using deep learning”. In: *eLife* 8 (2019), e47994. doi: 10.7554/elife.47994.
- 722 [23] Mayank Kabra et al. “JAABA: interactive machine learning for automatic annotation of
723 animal behavior”. In: *Nature Methods* 10.1 (2013), pp. 64–67. issn: 1548-7091. doi: 10.1038/
724 nmeth.2281.
- 725 [24] Caleb Weinreb et al. “Keypoint-MoSeq: parsing behavior by linking point tracking to
726 pose dynamics”. In: *Nature Methods* 21.7 (2024), pp. 1329–1339. issn: 1548-7091. doi:
727 10.1038/s41592-024-02318-2.
- 728 [25] Alexander I. Hsu and Eric A. Yttri. “B-SOIID, an open-source unsupervised algorithm for
729 identification and fast prediction of behaviors”. In: *Nature Communications* 12.1 (2021),
730 p. 5188. doi: 10.1038/s41467-021-25420-x.
- 731 [26] Kristin Branson et al. “High-throughput ethomics in large groups of *Drosophila*”. In: *Nature
732 Methods* 6.6 (2009), pp. 451–457. issn: 1548-7091. doi: 10.1038/nmeth.1328.
- 733 [27] Duncan S. Mearns et al. “Deconstructing Hunting Behavior Reveals a Tightly Coupled
734 Stimulus-Response Loop”. In: *Current Biology* 30.1 (2020), 54–69.e9. issn: 0960-9822. doi:
735 10.1016/j.cub.2019.11.022.
- 736 [28] Rebecca Z. Weber et al. “Deep learning-based behavioral profiling of rodent stroke recovery”.
737 In: *BMC Biology* 20.1 (2022), p. 232. doi: 10.1186/s12915-022-01434-9.
- 738 [29] Mido Assran et al. “V-JEPA 2: Self-Supervised Video Models Enable Understanding,
739 Prediction and Planning”. In: *arXiv* (2025). doi: 10.48550/arxiv.2506.09985. eprint:
740 2506.09985.

- 741 [30] Bin Lin et al. “Video-LLaVA: Learning Unified Visual Representation by Alignment Before
742 Projection”. In: *Proceedings of the 2024 Conference on Empirical Methods in Natural*
743 *Language Processing* (2024), pp. 5971–5984. doi: 10.18653/v1/2024.emnlp-main.342.
- 744 [31] Olivier Friard and Marco Gamba. “BORIS: a free, versatile open-source event-logging
745 software for video/audio coding and live observations”. In: *Methods in Ecology and Evolution*
746 7.11 (2016), pp. 1325–1330. issn: 2041-210X. doi: 10.1111/2041-210x.12584.
- 747 [32] Lucas P. J. J. Noldus, Andrew J. Spink, and Ruud A. J. Tegelenbosch. “EthoVision: A
748 versatile video tracking system for automation of behavioral experiments”. In: *Behavior*
749 *Research Methods, Instruments, & Computers* 33.3 (2001), pp. 398–414. issn: 0743-3808.
750 doi: 10.3758/bf03195394.
- 751 [33] Annemarie van der Marel et al. “A comparison of low-cost behavioral observation software
752 applications for handheld computers and recommendations for use”. In: *Ethology* 128.3
753 (2022), pp. 275–284. issn: 0179-1613. doi: 10.1111/eth.13251.
- 754 [34] Rishi Bommasani et al. “On the Opportunities and Risks of Foundation Models”. In: *arXiv*
755 (2021). doi: 10.48550/arxiv.2108.07258. eprint: 2108.07258.
- 756 [35] Jennifer J Sun et al. “The Multi-Agent Behavior Dataset: Mouse Dyadic Social Interactions”.
757 In: *arXiv* (2021). doi: 10.48550/arxiv.2104.02710. eprint: 2104.02710.
- 758 [36] Jennifer J Sun et al. “MABe22: A Multi-Species Multi-Task Benchmark for Learned
759 Representations of Behavior”. In: *arXiv* (2022). doi: 10.48550/arxiv.2207.10553. eprint:
760 2207.10553.
- 761 [37] Long Zhao et al. “VideoPrism: A Foundational Visual Encoder for Video Understanding”. In:
762 *arXiv* (2024). doi: 10.48550/arxiv.2402.13217. eprint: 2402.13217.
- 763 [38] Greg J. Stephens, Matthew Bueno de Mesquita, William S. Ryu, and William Bialek.
764 “Emergence of long timescales and stereotyped behaviors in *Caenorhabditis elegans*”. In:
765 *Proceedings of the National Academy of Sciences* 108.18 (2011), pp. 7286–7289. issn:
766 0027-8424. doi: 10.1073/pnas.1007868108.

- 767 [39] Neil A. Croll. “Components and patterns in the behaviour of the nematode *Caenorhabditis*
768 *elegans*”. In: *Journal of Zoology* 176.2 (1975), pp. 159–176. ISSN: 0952-8369. doi: 10.1111/j.
769 1469-7998.1975.tb03191.x.
- 770 [40] Benjamin Koger et al. “Quantifying the movement, behaviour and environmental context of
771 group-living animals using drones and computer vision”. In: *Journal of Animal Ecology* 92.7
772 (2023), pp. 1357–1371. ISSN: 0021-8790. doi: 10.1111/1365-2656.13904.
- 773 [41] Otto Brookes et al. “PanAf20K: A Large Video Dataset for Wild Ape Detection and Behaviour
774 Recognition”. In: *International Journal of Computer Vision* 132.8 (2024), pp. 3086–3102.
775 ISSN: 0920-5691. doi: 10.1007/s11263-024-02003-z.
- 776 [42] Vikram Chandra, Asaf Gal, and Daniel J. C. Kronauer. “Colony expansions underlie
777 the evolution of army ant mass raiding”. In: *Proceedings of the National Academy of
778 Sciences of the United States of America* 118.22 (2021), e2026534118. ISSN: 0027-8424. doi:
779 10.1073/pnas.2026534118.
- 780 [43] Alfonso Pérez-Escudero et al. “idTracker: tracking individuals in a group by automatic
781 identification of unmarked animals.” In: *Nature methods* 11.7 (2013), pp. 743–8. ISSN:
782 1548-7091. doi: 10.1038/nmeth.2994.
- 783 [44] Felix B Mueller, Timo Lueddecke, Richard Vogg, and Alexander S Ecker. “Domain-Adaptive
784 Pretraining Improves Primate Behavior Recognition”. In: *arXiv* (2025). doi: 10.48550/arxiv.
785 2509.12193. eprint: 2509.12193.
- 786 [45] S Brenner. “The Genetics of *Caenorhabditis elegans*”. In: *Genetics* 77.1 (1974), pp. 71–94.
787 ISSN: 0016-6731. doi: 10.1093/genetics/77.1.71.
- 788 [46] Aylesse Sordillo and Cornelia I Bargmann. “Behavioral control by depolarized and hyperpo-
789 larized states of an integrating neuron”. In: *eLife* 10 (2021), e67723. doi: 10.7554/elife.67723.
- 790 [47] Alejandro López-Cruz et al. “Parallel Multimodal Circuits Control an Innate Foraging
791 Behavior”. In: *Neuron* 102.2 (2019), 407–419.e8. ISSN: 0896-6273. doi: 10.1016/j.neuron.
792 2019.01.053.

- 793 [48] Kuang-Man Huang, Pamela Cosman, and William R. Schafer. “Machine vision based detection
794 of omega bends and reversals in *C. elegans*”. In: *Journal of Neuroscience Methods* 158.2
795 (2006), pp. 323–336. issn: 0165-0270. doi: 10.1016/j.jneumeth.2006.06.007.
- 796 [49] Eviatar Yemini et al. “A database of *Caenorhabditis elegans* behavioral phenotypes”. In:
797 *Nature Methods* 10.9 (2013), pp. 877–879. issn: 1548-7091. doi: 10.1038/nmeth.2560.
- 798 [50] Maria Teresa Lino et al. “Annotating Multi-media/Multi-modal Resources with ELAN”. In:
799 *Proceedings of the Fourth International Conference on Language Resources and Evaluation*
800 (*LREC’04*). Lisbon, 2004.
- 801 [51] Yi Wang et al. “InternVideo2: Scaling Foundation Models for Multimodal Video Understand-
802 ing”. In: *arXiv* (2024). doi: 10.48550/arxiv.2403.15377. eprint: 2403.15377.
- 803 [52] Loubna Ben Allal et al. “SmolLM2: When Smol Goes Big – Data-Centric Training of a Small
804 Language Model”. In: *arXiv* (2025). doi: 10.48550/arxiv.2502.02737. eprint: 2502.02737.
- 805 [53] Yingwei Li, Yi Li, and Nuno Vasconcelos. “RESOUND: Towards Action Recognition Without
806 Representation Bias”. In: *Lecture Notes in Computer Science* (2018), pp. 520–535. issn:
807 0302-9743. doi: 10.1007/978-3-030-01231-1_32.
- 808 [54] Ashish Vaswani et al. “Attention Is All You Need”. In: *arXiv* (2017). doi: 10.48550/arxiv.
809 1706.03762. eprint: 1706.03762.
- 810 [55] Sergey Ioffe and Christian Szegedy. “Batch Normalization: Accelerating Deep Network
811 Training by Reducing Internal Covariate Shift”. In: *arXiv* (2015). doi: 10.48550/arxiv.1502.
812 03167. eprint: 1502.03167.
- 813 [56] Geoffrey E Hinton et al. “Improving neural networks by preventing co-adaptation of feature
814 detectors”. In: *arXiv* (2012). doi: 10.48550/arxiv.1207.0580. eprint: 1207.0580.
- 815 [57] Samuel G. Müller and Frank Hutter. “TrivialAugment: Tuning-free Yet State-of-the-Art Data
816 Augmentation”. In: *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*
817 00 (2021), pp. 754–762. doi: 10.1109/iccv48922.2021.00081.

- 818 [58] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. “mixup: Beyond
 819 Empirical Risk Minimization”. In: *arXiv* (2017). doi: 10.48550/arxiv.1710.09412. eprint:
 820 1710.09412.
- 821 [59] Diederik P Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In: *arXiv*
 822 (2014). doi: 10.48550/arxiv.1412.6980. eprint: 1412.6980.
- 823 [60] Daniel Morales-Brottons, Thijs Vogels, and Hadrien Hendrikx. “Exponential Moving Average
 824 of Weights in Deep Learning: Dynamics and Benefits”. In: *arXiv* (2024). doi: 10.48550/arxiv.
 825 2411.18704. eprint: 2411.18704.
- 826 [61] Thomas Wolf et al. “HuggingFace’s Transformers: State-of-the-art Natural Language Pro-
 827 cessing”. In: *arXiv* (2019). doi: 10.48550/arxiv.1910.03771. eprint: 1910.03771.
- 828 [62] Adam Paszke et al. “PyTorch: An Imperative Style, High-Performance Deep Learning
 829 Library”. In: *arXiv* (2019). doi: 10.48550/arxiv.1912.01703. eprint: 1912.01703.

830 **Extended Tables**

Dataset	mAP	F1	Precision	Recall
CalMS21	0.945	0.893	0.89	0.895
Ant colonies	0.978	0.928	0.936	0.920
C. elegans	0.932	0.863	0.820	0.921
Adult-larva ants	0.925	0.813	0.728	0.930
MaBE	0.970	0.920	0.951	0.894
Zebras	0.853	0.785	0.809	0.768
PanAf500	0.657	0.578	0.667	0.527

Table 1: **Performance metrics across all evaluated datasets.** All experiments were conducted using the default configuration.

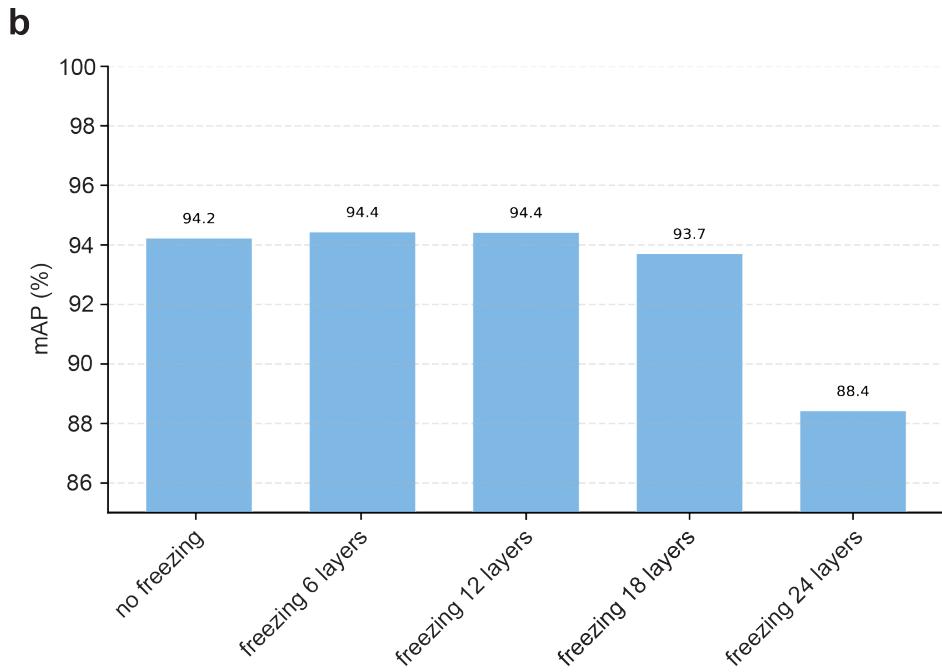
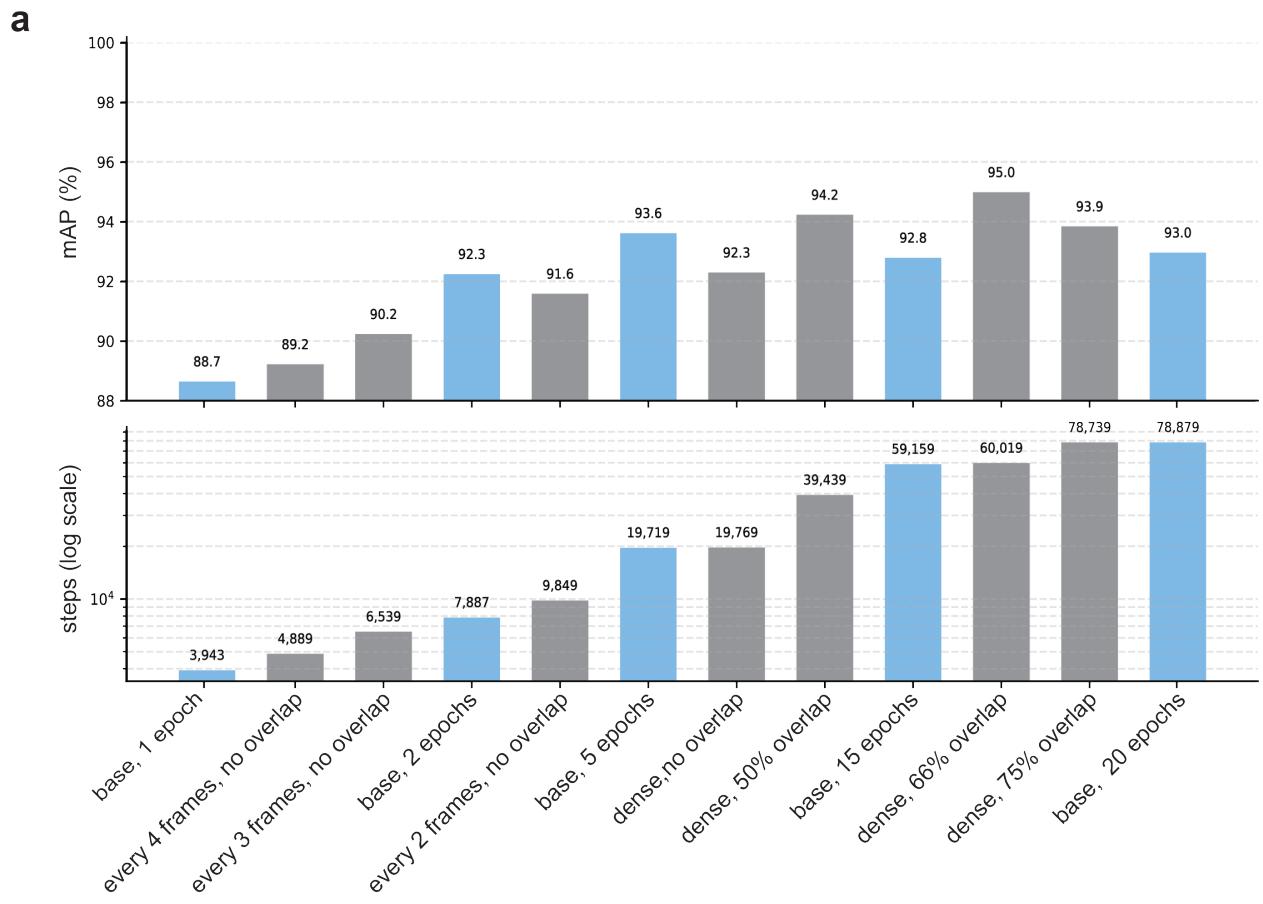
Dataset	Classes	Train videos	Train frames	Test frames	Test videos	Multilabel
CalMS21	4	70	507738	262107	19	False
Ant colonies	2	7	126000	36000	2	False
C. elegans	5	29	92560	35674	8	False
Adult-larva ants	3	7	126000	54000	3	False
MaBE	6	897	807300	202500	225	True
Zebras	3	35	2487594	424550	9	False
PanAf500	9	400	143959	35998	100	True

Table 2: Summary of datasets used in our experiments

Dataset	Class distribution (%)
CalMS21	<i>attack: 3.5%, investigate: 27.0%, mount: 7.9%, other: 61.7%</i>
Ant colonies	<i>other: 91.7%, raiding: 8.3%</i>
C. elegans	<i>other: 0.1%, forward: 84.6%, reverse: 7.3%, turn: 2.3%, pause: 5.8%</i>
Adult-larva ants	<i>other: 88.9%, self: 1.4%, larvae: 9.7%</i>
MaBE	<i>behavior 1: 21.1%, behavior 2: 15.5%, behavior 3: 21.3%, behavior 4: 7.8%, behavior 5: 24.6%, behavior 6: 9.6%</i>
Zebras	<i>other: 78.1%, out of sight: 2.4%, vigilant: 19.5%</i>
PanAf500	<i>camera interaction: 1.1%, climbing down: 0.8%, climbing up: 2.2%, hanging: 4.8%, running: 1.0%, sitting: 32.4%, sitting on back: 1.4%, standing: 21.6%, walking: 27.1%</i>

Table 3: Class frequency distribution (percentage of labeled frames across train and test splits) for each dataset.

831 Extended Figures



Extended Figure 1. Effect of temporal sampling and backbone freezing on FERAL performance. (a) Model quality as a function of number of training steps and overlapping-chunk configuration on the CalMS21 dataset. (b) Evaluation of layer-freezing strategies during fine-tuning.

832 **Supplementary Videos**

833 **Supplementary Video 1 | FERAL on CalMS21 mouse social interactions.** Representative
834 example from the CalMS21 dataset showing raw videos of resident–intruder mouse interactions
835 alongside frame-level FERAL predictions. [Link](#)

836 **Supplementary Video 2 | FERAL segmentation of *C. elegans* locomotor states.** Example
837 recordings of freely moving *C. elegans* with FERAL predictions overlaid on cropped worm-centered
838 views. [Link](#)

839 **Supplementary Video 3 | Adult–larva grooming behavior in clonal raider ants.** Example dyadic
840 interactions between adult clonal raider ants and larvae, with FERAL predictions overlaid on the
841 raw videos. [Link](#)

842 **Supplementary Video 4 | FERAL on PanAf500 camera-trap recordings of wild apes.** Sample
843 clip from the PanAf500 dataset showing a wild gorilla recorded by a camera trap in its natural
844 habitat. FERAL’s multi-label predictions are displayed frame by frame. [Link](#)