

1	63
2	64
3	65
4	66
5	67
6	68
7	69
8	70
9	71
10	72
11	73
12	74
13	75
14	76
15	77
16	78
17	79
18	80
19	81
20	82
21	83
22	Nonparametric inference of interaction laws in systems of agents from trajectory data
23	84
24	Fei Lu, Ming Zhong, Sui Tang, Mauro Maggioni
25	85
26	Mauro Maggioni.
27	E-mail: mauromaggionijhu@icloud.com
28	86
29	87
30	This PDF file includes:
31	93
32	Supplementary text
33	94
34	Figs. S1 to S17
35	95
36	Tables S1 to S22
37	96
38	References for SI reference citations
39	97
40	98
41	99
42	100
43	101
44	102
45	103
46	104
47	105
48	106
49	107
50	108
51	109
52	110
53	111
54	112
55	113
56	114
57	115
58	116
59	117
60	118
61	119
62	120
	121
	122
	123
	124

125	<b>Supporting Information Text</b>	187
126		188
127	<b>1. Learning Theory</b>	189
128	Consider the problem of estimating the interaction kernel $\phi : \mathbb{R}_+ \rightarrow \mathbb{R}$ of the dynamical system as follows	190
129		191
130		192
131	$\dot{\mathbf{x}}_i(t) = \frac{1}{N} \sum_{i'=1}^N \phi(\ \mathbf{x}_{i'}(t) - \mathbf{x}_i(t)\ ) (\mathbf{x}_{i'}(t) - \mathbf{x}_i(t)),$	[1] 193
132		194
133		195
134	from observations of discrete-time trajectories and derivatives, $\{\mathbf{X}^m(t_l)\}$ and $\{\dot{\mathbf{X}}^m(t_l)\}$ with $0 = t_1 < \dots < t_L = T$ and	196
135	$m = 1, \dots, M$ . We let $\mathbf{X} := (\mathbf{x}_i)_{i=1}^N \in \mathbb{R}^{dN}$ be the state space variable. The initial conditions $\mathbf{X}_0^m := \mathbf{X}^m(0)$ are sampled	197
136	independently from a probability measure $\mu_0$ on $\mathbb{R}^{dN}$ .	198
137	Such a system can also be described as the gradient flow $\dot{\mathbf{X}} = \mathbf{f}_\phi(\mathbf{X}) = \nabla \mathcal{U}(\mathbf{X})$ of the potential energy $\mathcal{U}(\mathbf{X}) =$	199
138	$\frac{1}{2N} \sum_{i,i'} \Phi(\ \mathbf{x}_i - \mathbf{x}_{i'}\ )$ , with the function $\Phi : \mathbb{R}_+ \rightarrow \mathbb{R}$ satisfying $\Phi'(r) = \phi(r)r$ . Therefore, the estimation of $\phi$ is equivalent to	200
139	the estimation of $\Phi'$ . As we will see later, the function $\phi(\cdot)$ appears naturally in assessing the quality of approximation of	201
140	estimators of $\phi$ , the fundamental reason being the relationship with the potential involving $\Phi$ .	202
141	We restrict our attention to kernels in the <i>admissible set</i>	203
142		204
143		205
144	$\mathcal{K}_{R,S} := \{\phi \in W^{1,\infty} : \text{supp}(\phi) \in [0, R], \sup_{r \in [0, R]} [ \phi(r)  +  \phi'(r) ] \leq S\}$	[2] 206
145		207
146	for some $R, S > 0$ . The boundedness of $\phi$ and its derivative ensures the existence and uniqueness of a global solution to initial	208
147	value problems of the first order system Eq. (1), and the continuous dependence of the solution on the initial condition. The	209
148	restriction $\text{supp}(\phi) \subset [0, R]$ represents the finite range of interaction between particles, and this restriction may be replaced by	210
149	functions with unbounded support but with a suitable decay on $\mathbb{R}_+$ .	211
150		212
151	We shall construct an error functional based on the special structure of the dynamical system $\dot{\mathbf{X}} = \mathbf{f}_\phi(\mathbf{X})$ , taking advance	213
152	of the form of the dependency of the right-hand side $\mathbf{f}_\phi$ on the interaction kernel $\phi$ . This learning procedure deviates from	214
153	standard regression in two aspects: (i) the values of the interaction kernel are not observed, and cannot be explicitly estimated	215
154	from the observations of the state variables; (ii) the observations of the independent variable of the interaction kernel, given by	216
155	the pairwise distance between the agents, though abundant, are not independent and may be redundant.	217
156	We would also like to stress the importance of using a carefully chosen measure on the pairwise distance space, so as to	218
157	account for both the randomness from the initial conditions and the evolution of the dynamical system, and to reflect the	219
158	(relative) abundances of pairwise distances. Our analysis shows that the expectation of the empirical measure of the pairwise	220
159	distances is a natural choice, and it is closely related to the coercivity condition, the other fundamental ingredient which	221
160	ensures learnability and convergence of the estimators.	222
161		223
162	<b>A. The Error functional and estimators.</b> Given the structure of the first order system Eq. (1), we consider the error functional	224
163		225
164		226
165	$\mathcal{E}_{L,M}(\varphi) := \frac{1}{MN} \sum_{l,m,i=1}^{L,M,N} w_l \ \dot{\mathbf{x}}_i^m(t_l) - \mathbf{f}_\varphi(\mathbf{x}^m(t_l))_i\ ^2,$	[3] 227
166		228
167		229
168	where $\{w_l\}_{l=1}^L$ is a normalized set of weights ( $w_l > 0$ and $\sum_{l=1}^L w_l = 1$ ), and define an estimator	230
169		231
170		232
171	$\hat{\phi}_{L,M,\mathcal{H}} := \arg \min_{\varphi \in \mathcal{H}} \mathcal{E}_{L,M}(\varphi),$	[4] 233
172		234
173		235
174	where $\mathcal{H}$ is a suitable class of functions that will be referred as hypothesis space. Natural choices of weights $\{w_l\}$ may be	236
175	chosen to be all equal to $1/L$ , as in the case of equi-spaced $t_l$ 's, which is what we considered throughout the paper, and is	237
176	consistent with the definition of $\rho_T^L$ and its use in measuring the performance of the estimator in $L^2(\rho_T^L)$ . However, if one	238
177	wished to measure the performance in a different $L^2$ space, one could choose the weights differently. A distinguished choice	239
178	would be $L^2(\rho_{\text{Lebesgue}})$ , in which case one may choose $w_l = 1/(t_{l+1} - t_l)$ , for $l = 1, \dots, L - 1$ (and change all the summations	240
179	involving $l$ to stop at $L - 1$ instead of $L$ ). Other choices of weights corresponding to other quadrature rules are also be possible.	241
180	Note that the error functional is quadratic in $\varphi$ and bounded below by 0, therefore the minimizer exists for any finite	242
181	dimensional convex hypothesis spaces $\mathcal{H}$ . We can always truncate this minimizer so that it is bounded above by $S$ , the upper	243
182	bound of the functions in the admissible set $\mathcal{K}_{R,S}$ , and this truncated estimator behaves similarly to the estimator obtained by	244
183	assuming that the functions in $\mathcal{H}$ are uniformly bounded. In fact, such truncation can only reduce the error. Hence, without	245
184	loss of generality, we assume $\mathcal{H}$ to be a compact set in the $L^\infty$ norm.	246
185	Our objectives are measuring the quality of approximation of the estimator and finding the hypothesis spaces for which the	247
186	optimal rate of convergence of $\hat{\phi}$ to the true interaction kernel $\phi$ is achieved.	248

249	<b>B. Measures on the pairwise distance space.</b> We introduce a probability measure on $\mathbb{R}_+$ , to define a suitable function space	311
250	that contains all the estimators and the true interaction kernel, and to provide a norm to assess the accuracy of the estimators.	312
251	We let	313
252	$\mathbf{r}_{ii'}(t) = \mathbf{x}_{i'}(t) - \mathbf{x}_i(t)$ , and $r_{ii'}(t) = \ \mathbf{x}_{i'}(t) - \mathbf{x}_i(t)\ $ .	314
253	Note that the independent variable of the interaction kernel is the pairwise distances $r_{ii'}^m(t)$ , which can be computed from the	315
254	observed trajectories. It is natural to start from the empirical measure of pairwise distances	316
255		317
256	$\rho_T^{L,M}(r) = \frac{1}{\binom{N}{2} LM} \sum_{l,m=1}^{L,M} \sum_{i,i'=1, i < i'}^N \delta_{r_{ii'}^m(t_l)}(r)$ ,	[5] 318
257		319
258		320
259	which tends, as $M \rightarrow \infty$ , using the law of large numbers, to $\rho_T^L$ defined in (5) in the main text. When trajectories are observed	321
260	continuously in time, the counterpart of $\rho_T^L$ is the measure defined in (5). We now establish basic properties of these measures:	322
261	<b>Lemma 1.1.</b> <i>For each <math>\phi \in \mathcal{K}_{R,S}</math> defined in Eq. (2), the measures <math>\rho_T^L</math> and <math>\rho_T</math> defined in (5) and (4) in the main text are Borel</i>	323
262	<i>probability measures on <math>\mathbb{R}_+</math>. They are absolutely continuous with respect to the Lebesgue measure provided that <math>\mu_0</math> is absolutely</i>	324
263	<i>continuous with respect to the Lebesgue measure on <math>\mathbb{R}^{dN}</math>.</i>	325
264		326
265	<b>C. Learnability: the coercivity condition.</b> A fundamental question is the learnability of the true interaction kernel, i.e. the	327
266	well-posedness of the inverse problem of kernel learning. Since the estimators $\hat{\phi}_{L,M,\mathcal{H}}$ always exists for suitably chosen hypothesis	328
267	spaces $\mathcal{H}$ (e.g. compact sets), learnability is equivalent to the convergence of the estimator $\hat{\phi}_{L,M,\mathcal{H}}$ to the true kernel $\phi$ as the	329
268	sample size increases (i.e. $M \rightarrow \infty$ ) and as the hypothesis space grows. To ensure such a convergence, one would naturally	330
269	wish: (i) that the true kernel $\phi$ is the unique minimizer of the expectation of the error functional (by the law of large numbers)	331
270		332
271	$\mathcal{E}_{L,\infty}(\varphi) := \lim_{M \rightarrow \infty} \mathcal{E}_{L,M}(\varphi) = \frac{1}{LN} \sum_{l,i=1}^{L,N} \mathbb{E} \left[ \left\  \frac{1}{N} \sum_{i'=1}^N (\varphi - \phi)(r_{ii'}(t_l)) \mathbf{r}_{ii'}(t_l) \right\ ^2 \right]$ ,	[6] 333
272		334
273		335
274	(ii) that the error of the estimator, in terms of a metric based on the $L^2(\rho_T^L)$ norm, can be controlled by the discrepancy	336
275	between the empirical error functional and its limit.	337
276	Note that $\mathcal{E}_{L,\infty}(\varphi) \geq 0$ for any $\varphi$ and that $\mathcal{E}_{L,\infty}(\phi) = 0$ . Furthermore, Eq. (6) reveals that $\mathcal{E}_{L,\infty}(\varphi)$ is a quadratic functional	338
277	of $\varphi - \phi$ , and we have, by Jensen's inequality,	339
278		340
279	$\mathcal{E}_{L,\infty}(\varphi) \leq \frac{(N-1)^2}{N^2} \ \varphi(\cdot) - \phi(\cdot)\ _{L^2(\rho_T^L)}^2$ .	341
280		342
281	This inequality suggests the above weighted $L^2(\rho_T^L)$ norm as a metric on the error of the estimator that we wish to be controlled.	343
282	Therefore, as long we as can bound the limit error functional from below by $\ \varphi(\cdot) - \phi(\cdot)\ _{L^2(\rho_T^L)}^2$ , we can conclude that $\phi$ is	344
283	the unique minimizer of $\mathcal{E}_{L,\infty}(\cdot)$ and that the estimators converge to $\phi$ . This suggests the following coercivity condition:	345
284	<b>Definition 1.1</b> (Coercivity condition). <i>We say that the dynamical system defined in Eq. (1) together with the probability</i>	346
285	<i>measure <math>\mu_0</math> on <math>\mathbb{R}^{dN}</math>, satisfies the coercivity condition on <math>\mathcal{H}</math> with a constant <math>c_{L,N,\mathcal{H}} &gt; 0</math>, if</i>	347
286		348
287	$c_{L,N,\mathcal{H}} \ \varphi(\cdot) - \phi(\cdot)\ _{L^2(\rho_T^L)}^2 \leq \frac{1}{NL} \sum_{i,l=1}^{L,N} \mathbb{E} \left[ \left\  \frac{1}{N} \sum_{i'=1}^N \varphi(r_{ii'}(t_l)) \mathbf{r}_{ii'}(t_l) \right\ ^2 \right]$	[7] 349
288		350
289		351
290	for all $\varphi \in \mathcal{H}$ such that $\varphi(\cdot) \in L^2(\rho_T^L)$ , with the measure $\rho_T^L$ defined in (4) in the main text, and the expectation being with	352
291	respect to initial conditions distributed according to $\mu_0$ .	353
292	The above inequality is called a coercivity condition because that it implies coercivity of the bilinear functional $\langle\langle \cdot, \cdot \rangle\rangle$ on	354
293	$L^2(\mathbb{R}_+, \rho_T^L)$ ,	355
294		356
295	$\langle\langle \varphi_1, \varphi_2 \rangle\rangle := \frac{1}{LN} \sum_{l,i=1}^{L,N} \mathbb{E} \left[ \left\langle \frac{1}{N} \sum_{j=1}^N \varphi_1(r_{ji}(t_l)) \mathbf{r}_{ij}(t_l), \frac{1}{N} \sum_{j=1}^N \varphi_2(r_{ji}(t_l)) \mathbf{r}_{ij}(t_l) \right\rangle \right]$ ,	[8] 357
296		358
297		359
298	as Eq. (7) may be rewritten as	360
299	$c_{L,N,\mathcal{H}} \ \varphi(\cdot) - \phi(\cdot)\ _{L^2(\mathbb{R}_+, \rho_T^L)}^2 \leq \langle\langle \varphi, \varphi \rangle\rangle$ .	361
300		362
301	The coercivity condition plays a key role in the learning of the interaction kernel. It ensures learnability by ensuring the	363
302	uniqueness of minimizer of the expectation of the error functional, and by guaranteeing convergence of estimators through a	364
303	control of the error of the estimator on every compact convex hypothesis space $\mathcal{H}$ in $L^2(\rho_T^L)$ . To see this, apply the coercivity	365
304	inequality to $\varphi - \phi$ , to obtain	366
305	$c_{L,N,\mathcal{H}} \ \varphi(\cdot) - \phi(\cdot)\ _{L^2(\mathbb{R}_+, \rho_T^L)}^2 \leq \mathcal{E}_{L,\infty}(\varphi)$ .	[9] 367
306		368
307	From the facts that $\mathcal{E}_{L,\infty}(\varphi) \geq 0$ for any $\varphi$ and that $\mathcal{E}_{L,\infty}(\phi) = 0$ , we can conclude that the true kernel $\phi$ is the unique minimizer	369
308	of the $\mathcal{E}_{L,\infty}(\varphi)$ . Furthermore, the coercivity condition enables us to control the error of the estimator, on every compact convex	370
309	hypothesis space in $L^2(\rho_T^L)$ , by the discrepancy of the error functional (see Proposition 1.3), therefore guaranteeing convergence	371
310	of the estimator.	372

373 **Theorem 1.2.** Let  $\mathcal{H}_n$  be a sequence of compact convex subsets of  $L^\infty([0, R])$  such that 435  
 374  $\inf_{\varphi \in \mathcal{H}_n} \|\varphi(\cdot) - \phi(\cdot)\|_{L^2(\rho_T^L)} \rightarrow 0$  436  
 375 as  $n \rightarrow \infty$ . Assume that the coercivity condition holds on  $\cup_{n=1}^\infty \mathcal{H}_n$ . Then the estimator  $\hat{\phi}_{L,M,\mathcal{H}_n}$  defined in Eq. (4) converges 437  
 376 to the true kernel in  $L^2(\rho_T^L)$  almost surely as  $n, M$  approaches infinity, i.e. 438  
 377  $\lim_{n \rightarrow \infty} \lim_{M \rightarrow \infty} \|\hat{\phi}_{L,M,\mathcal{H}_n}(\cdot) - \phi(\cdot)\|_{L^2(\rho_T^L)} = 0$ , almost surely. 439  
 378  
 379 The above theorem follows from the next proposition. 440  
 380  
 381 **Proposition 1.3.** Let  $\mathcal{H}$  be a compact convex subset of  $L^2(\rho_T^L)$  and assume the coercivity condition holds true on  $\mathcal{H}$ . Then the 441  
 382 functional  $\mathcal{E}_{L,\infty}$  defined in Eq. (6) admits a unique minimizer 442  
 383  $\hat{\phi}_{L,\infty,\mathcal{H}} = \arg \min_{\varphi \in \mathcal{H}} \mathcal{E}_{L,\infty}(\varphi)$ , 443  
 384 in  $L^2(\rho_T^L)$ . Furthermore, for all  $\varphi \in \mathcal{H}$  444  
 385  $\mathcal{E}_{L,\infty}(\varphi) - \mathcal{E}_{L,\infty}(\hat{\phi}_{L,\infty,\mathcal{H}}) \geq c_{L,N,\mathcal{H}} \|\varphi(\cdot) - \hat{\phi}_{L,\infty,\mathcal{H}}(\cdot)\|_{L^2(\rho_T^L)}^2$ . 445  
 386  
 387 **D. Optimal rate of convergence of the estimator.** We now turn to the rate of convergence of the estimator. 446  
 388 **Theorem 1.4.** Let the true kernel  $\phi \in \mathcal{K}_{R,S}$ , and let  $\mathcal{H} \subset L^\infty([0, R])$  be compact convex and bounded above by  $S_0 \geq S$ . Assume 447  
 389 that the coercivity condition in Eq. (7) holds. Then for any  $\epsilon > 0$ , we have 448  
 390  $c_{L,N,\mathcal{H}} \|\hat{\phi}_{L,M,\mathcal{H}}(\cdot) - \phi(\cdot)\|_{L^2(\rho_T^L)}^2 \leq 2 \inf_{\varphi \in \mathcal{H}} \|\varphi(\cdot) - \phi(\cdot)\|_{L^\infty([0, R])}^2 + 2\epsilon$  449  
 391 with probability at least  $1 - \delta$ , provided that 450  
 392  $M \geq \frac{1152S_0^2R^2}{c_{L,N,\mathcal{H}}\epsilon} \left( \log(\mathcal{N}(\mathcal{H}, \frac{\epsilon}{48S_0R^2})) + \log(\frac{1}{\delta}) \right)$ , 451  
 393 where  $\mathcal{N}(\mathcal{H}, \eta)$  is the  $\eta$ -covering number of  $\mathcal{H}$  under the  $\infty$ -norm. 452  
 394  
 395 We discuss first the implications of this theorem on the choice of hypothesis space in view of obtaining optimal rates of 453  
 396 convergence of our estimator. The proof of the theorem will be presented at the end of this section. In practice, given a 454 set of  $M$  trajectories, we would like to chose the best finite-dimensional hypothesis space  $\mathcal{H}$  to minimize the error of the 455 estimator. There are two competing issues. On one hand, we would like the hypothesis space  $\mathcal{H}$  to be large so that the 456 bias  $\inf_{\varphi \in \mathcal{H}} \|\varphi - \phi\|_{L^\infty([0, R])}^2$  is small. On the other hand, we would like to keep  $\mathcal{H}$  to be small so that the covering number 457  $\mathcal{N}(\mathcal{H}, \epsilon/48S_0R^2)$ , and therefore the variance of the estimator is small. This is the classical bias-variance trade-off in statistical 458 estimation. Inspired from approximation methods in regression (1–3), the following proposition quantifies the effect of 459 hypothesis spaces on the rate of convergence of the estimator. 460  
 397  
 398 **Proposition 1.5.** Assume that the coercivity condition holds with a constant  $c_{L,N,\mathcal{H}}$ , and recall  $\hat{\phi}_{L,M,\mathcal{H}}$  defined in Eq. (4) is a 461  
 399 minimizer of the empirical error functional over a hypothesis space  $\mathcal{H}$ . 462  
 400 (a) For  $\mathcal{H} = \mathcal{K}_{R,S}$ , there exists a constant  $C = C(S, R)$  such that 463  
 401  
 402 (b) Assume that  $\mathcal{H}_n$  is a sequence of finite dimensional spaces of  $L^\infty([0, R])$  such that  $\dim(\mathcal{H}_n) \leq c_{\text{on}}$  and 464  
 403  $\inf_{\varphi \in \mathcal{H}_n} \|\varphi(\cdot) - \phi(\cdot)\|_{L^\infty([0, R])}^2 \leq c_1 n^{-s}$  465  
 404 for all  $n$  for some constants  $c_0, c_1, s > 0$ , then by choosing  $n = n_* := (M/\log M)^{\frac{1}{2s+1}}$ , we have 466  
 405  $\mathbb{E}[\|\hat{\phi}_{L,M,\mathcal{H}_{n_*}}(\cdot) - \phi(\cdot)\|_{L^2(\rho_T^L)}] \leq \frac{C}{c_{L,N,\mathcal{H}}} \left( \frac{\log M}{M} \right)^{\frac{s}{2s+1}}$ , 467  
 406 where  $C = C(c_0, c_1, R, S)$ . 468  
 407  
 408 It is interesting to compare this rate with those in the mean field regime, where the regime  $N \rightarrow \infty$  (with  $M = 1, L \rightarrow \infty$ ) 469  
 409 was studied: the rates implied by (4) they seem to be no better than  $N^{-1/d}$ , i.e. they are cursed by the dimension  $d$ , even if 470 the problem is fundamentally that of estimating a 1-dimensional function. It would be interesting to understand whether that 471 rate is optimal for this problem in the mean-field regime ( $N \rightarrow \infty$ ), or if in fact, the results in the present work lead to sharper, 472 dimension-independent bounds in the mean-field limit as well. 473  
 413 The proof of Thm. 1.4 is based on this technical Proposition: 474  
 414

497 **Proposition 1.6.** Assume the coercivity condition holds true and let  $\mathcal{H} \subset L^\infty([0, R])$  be compact convex, bounded above by  $S_0$ .  
 498 Let  
 499 
$$\mathcal{D}_{L,\infty,\mathcal{H}}(\varphi) := \mathcal{E}_{L,\infty}(\varphi) - \mathcal{E}_{L,\infty}(\widehat{\phi}_{L,\infty,\mathcal{H}}) , \quad \mathcal{D}_{L,M,\mathcal{H}}(\varphi) := \mathcal{E}_{L,M}(\varphi) - \mathcal{E}_{L,M}(\widehat{\phi}_{L,\infty,\mathcal{H}}),$$
  
 500 where  $\widehat{\phi}_{L,\infty,\mathcal{H}}$  is the minimizer of  $\mathcal{E}_{L,\infty}(\cdot)$  over  $\mathcal{H}$ . Then for all  $\epsilon > 0$  and  $0 < \alpha < 1$ , we have  
 501 
$$P\left\{ \sup_{\varphi \in \mathcal{H}} \frac{\mathcal{D}_{L,\infty,\mathcal{H}}(\varphi) - \mathcal{D}_{L,M,\mathcal{H}}(\varphi)}{\mathcal{D}_{L,\infty,\mathcal{H}}(\varphi) + \epsilon} \geq 3\alpha \right\} \leq \mathcal{N}(\mathcal{H}, C_1 \alpha \epsilon) e^{-C_2 \alpha^2 M \epsilon}$$
  
 502 where  $C_1 = \frac{1}{8S_0 R^2}$  and  $C_2 = \frac{-c_{L,N,\mathcal{H}}}{32S_0^2 R^2}$ .  
 503 **Proof of the Theorem 1.4 .** Put  $\alpha = \frac{1}{6}$  in Proposition 1.6. We know that, with probability at least  
 504 
$$1 - \mathcal{N}\left(\mathcal{H}, \frac{\epsilon}{48S_0 R^2}\right) e^{-\frac{c_{L,N,\mathcal{H}} M \epsilon}{1152 S_0^2 R^2}},$$
  
 505 we have  
 506 and therefore, for all  $\varphi \in \mathcal{H}$ ,  
 507 
$$\frac{1}{2} \mathcal{D}_{L,\infty,\mathcal{H}}(\varphi) < \mathcal{D}_{L,M,\mathcal{H}}(\varphi) + \frac{1}{2} \epsilon.$$
  
 508 Taking  $\varphi = \widehat{\phi}_{L,M,\mathcal{H}}$ , we have  
 509 
$$\mathcal{D}_{L,\infty,\mathcal{H}}(\widehat{\phi}_{L,M,\mathcal{H}}) < 2\mathcal{D}_{L,M,\mathcal{H}}(\widehat{\phi}_{L,M,\mathcal{H}}) + \epsilon.$$
  
 510 But  $\mathcal{D}_{L,M,\mathcal{H}}(\widehat{\phi}_{L,M,\mathcal{H}}) = \mathcal{E}_{L,M}(\widehat{\phi}_{L,M,\mathcal{H}}) - \mathcal{E}_{L,M}(\widehat{\phi}_{L,\infty,\mathcal{H}}) \leq 0$  and hence by Proposition 1.3 we have  
 511 
$$c_{L,N,\mathcal{H}} \|\widehat{\phi}_{L,M,\mathcal{H}}(\cdot) - \widehat{\phi}_{L,\infty,\mathcal{H}}(\cdot)\|_{L^2(\rho_T^L)}^2 \leq \mathcal{D}_{L,\infty,\mathcal{H}}(\widehat{\phi}_{L,M,\mathcal{H}}) < \epsilon.$$
  
 512 Therefore,  
 513 
$$\begin{aligned} \|\widehat{\phi}_{L,M,\mathcal{H}}(\cdot) - \phi(\cdot)\|_{L^2(\rho_T^L)}^2 &\leq 2\|\widehat{\phi}_{L,M,\mathcal{H}}(\cdot) - \widehat{\phi}_{L,\infty,\mathcal{H}}(\cdot)\|_{L^2(\rho_T^L)}^2 + 2\|\widehat{\phi}_{L,\infty,\mathcal{H}}(\cdot) - \phi(\cdot)\|_{L^2(\rho_T^L)}^2 \\ &\leq \frac{2}{c_{L,N,\mathcal{H}}} (\epsilon + \inf_{\varphi \in \mathcal{H}} \|\varphi(\cdot) - \phi(\cdot)\|_\infty^2), \end{aligned}$$
  
 514 where the last inequality follows from the coercivity condition and by the definition of  $\widehat{\phi}_{L,\infty,\mathcal{H}}$ (see Eq. (10)). Given  $0 < \delta < 1$ ,  
 515 we see we need  $M$  large enough so that  
 516 
$$1 - \mathcal{N}\left(\mathcal{H}, \frac{\epsilon}{48S_0 R^2}\right) e^{-\frac{c_{L,N,\mathcal{H}} M \epsilon}{1152 S_0^2 R^2}} \geq 1 - \delta.$$
  
 517 The conclusion follows. □  
 518 **E. Trajectory-based Performance Measures.** After having established results on the convergence rate of our estimator, we turn  
 519 to control the accuracy of trajectories predicted when using the estimated interaction kernel, evolved from initial conditions  
 520 both in and outside of the training data. Trajectory-based measurements of accuracy are interesting because (a) they provide a  
 521 quantitative assessment on the quality of the approximated dynamics, (b) while the true interactions kernels are typically not  
 522 known, and so the accuracy of the estimated interaction kernel may not be evaluated, trajectories are known, and may be used  
 523 to perform model validation and cross-validation for parameter selection (if needed).  
 524 The next Proposition shows that the error in prediction is (i) bounded trajectory-wise by a continuous time version of the  
 525 error functional, and (ii) bounded in the mean squared sense by the mean squared error of the estimated interaction kernel.  
 526 **Proposition 1.7.** Let  $\widehat{\phi}$  be an estimator of the true interaction kernel  $\phi$ . Suppose that the function  $\widehat{\phi}(\|\cdot\|)$  is Lipschitz  
 527 continuous on  $\mathbb{R}^d$ , with Lipschitz constant  $C_{\text{Lip}}$ . Denote by  $\widehat{\mathbf{X}}(t)$  and  $\mathbf{X}(t)$  the solutions of the systems with interaction kernels  
 528  $\widehat{\phi}$  and  $\phi$  respectively, starting from the same initial condition. Then we have  
 529 
$$\sup_{t \in [0, T]} \|\widehat{\mathbf{X}}(t) - \mathbf{X}(t)\|^2 \leq 2T e^{8T^2 C_{\text{Lip}}^2} \int_0^T \|\dot{\mathbf{X}}(t) - \mathbf{f}_{\widehat{\phi}}(\mathbf{X}(t))\|^2 dt$$
  
 530 for each trajectory, and on average with respect to the initial distribution  $\mu_0$ ,  
 531 
$$\mathbb{E}_{\mu_0} [\sup_{t \in [0, T]} \|\widehat{\mathbf{X}}(t) - \mathbf{X}(t)\|^2] \leq C(T, C_{\text{Lip}}) \sqrt{N} \|\widehat{\phi}(\cdot) - \phi(\cdot)\|_{L^2(\rho_T)}^2$$
  
 532 for a constant  $C(T, C_{\text{Lip}})$ , where the measure  $\rho_T$  is as in Eq. (4) in the main text.

621	<b>2. Algorithm</b>	683
622	We start from describing the algorithm in its simplest form, for learning first order system with homogeneous agents; we then	684
623	move to first order systems with heterogeneous agents, and finish with the second order systems with heterogeneous agents.	685
624		686
625		687
626	<b>A. First Order Homogeneous Agent Systems.</b> Recall that we would like to estimate the interaction kernel $\phi$ of the $N$ -agent	688
627	system in Eq. (1) from $M$ independent trajectories $\{\mathbf{x}_i^m(t_l), \dot{\mathbf{x}}_i^m(t_l)\}_{i=1, l=1, m=1}^{N, L, M}$ with $t_l = \frac{lT}{L}$ . We obtain an estimator by	689
628	minimizing the discrete empirical error functional, over all $\varphi$ in a hypothesis space $\mathcal{H}_n$ ,	690
629		691
630	$\mathcal{E}_{L,M}(\varphi) = \frac{1}{LMN} \sum_{l,m,i=1}^{L,M,N} \left\  \dot{\mathbf{x}}_i^m(t_l) - \sum_{i'=1}^N \frac{1}{N} \varphi(r_{i,i'}^m(t_l)) \mathbf{r}_{i,i'}^m(t_l) \right\ ^2. \quad [14]$	692
631		693
632		694
633		695
634	When only the positions can be observed, we assume that $T/L$ is sufficiently small so that we can accurately approximate	696
635	the velocity $\dot{\mathbf{x}}_i^m(t_l)$ by finite differences, for example	697
636		698
637	$\dot{\mathbf{x}}_i^m(t_l) \approx \Delta \mathbf{x}_i^m(t_l) = \frac{\mathbf{x}_i^m(t_l) - \mathbf{x}_i^m(t_{l-1})}{t_l - t_{l-1}}, \quad \text{for } 1 \leq l \leq L,$	699
638		700
639		701
640	where we assumed $t_0$ is also observed. The error of the backward difference approximation is of order $O(T/L)$ , leading to a	702
641	$O(T/L)$ bias in the estimator. Therefore, for simplicity, we assume in the theoretical discussion that follows that the velocity	703
642	$\dot{\mathbf{x}}_i^m(t_l)$ is observed.	704
643		705
644	First, we set the hypothesis space $\mathcal{H}_n$ to be the span of $\{\psi_p\}_{p=1}^n$ , a set of linearly independent functions on $[0, R]$ . It is	706
645	natural to use an orthonormal basis of $\mathcal{H}_n$ in $L^2(\rho_L^T)$ for efficient computations. If the true interaction kernel is known to be	707
646	uniformly smooth, a global basis (e.g. Fourier) may be used. Since our admissible set is in $W^{1,\infty}$ , we shall use a local basis	708
647	consisting of piecewise polynomial functions on a partition of increasingly finer intervals. The partitions will be on the interval	709
648	$[R_{min}, R_{max}]$ , where $R_{min}$ and $R_{max}$ are minimal and maximal values of $r$ such that the empirical density $\rho_{L,M}^T(r)$ of the	710
649	pairwise distances $\{r_{i,i'}^m(t_l)\}$ is greater than a threshold.	711
650	Next, we minimize the empirical error functional over $\mathcal{H}_n$ to obtain an estimator. To simplify notation, for each $m$ , we	712
651	denote	713
652	$\mathbf{d}^m := (\dot{\mathbf{x}}_1^m(t_2), \dots, \dot{\mathbf{x}}_N^m(t_2); \dots; \dot{\mathbf{x}}_1^m(t_L) \dots \dot{\mathbf{x}}_N^m(t_L)) \quad [15]$	714
653		715
654	a column vector in $\mathbb{R}^{LNd}$ ; and denote	716
655		717
656	$\Psi_L^m(l, p) := \sum_{i'=1}^N \frac{1}{N} \psi_p(r_{i,i'}^m(t_l)) \mathbf{r}_{i,i'}^m(t_l) \in \mathbb{R}^d,$	718
657		719
658		720
659		721
660	for $1 \leq l \leq L$ , $1 \leq i \leq N$ and $1 \leq p \leq n$ , and refer it as the learning matrix $\Psi_L^m$ . Here and in what follows, the index $li$ denotes,	722
661	with some abuse of notation, the double-index $(l, i)$ mapped (in any fixed way) bijectively onto a one-dimensional array. Then	723
662	we can rewrite the empirical error functional as	724
663		725
664	$\mathcal{E}_{L,M}(\varphi) = \mathcal{E}_{L,M}(\mathbf{a}) = \frac{1}{LN M} \sum_{m=1}^M \ \mathbf{d}^m - \Psi_L^m \mathbf{a}\ _{\mathbb{R}^{LNd}}^2.$	726
665		727
666		728
667		729
668	Our estimator is the minimizer of $\mathcal{E}_{L,M}(\mathbf{a})$ over $\mathbb{R}^n$ . This is a Least Squares problem, and we solve for the minimizer from the	730
669	normal equations	731
670		732
671	$\underbrace{\frac{1}{M} \sum_{m=1}^M A_L^m \mathbf{a}}_{A_{L,M}} = \frac{1}{M} \sum_{m=1}^M b_L^m, \quad [16]$	733
672		734
673		735
674		736
675	where the trajectory-wise regression matrices are	737
676		738
677	$A_L^m := \frac{1}{LN} (\Psi_L^m)^T \Psi_L^m \quad \text{and} \quad b_L^m := \frac{1}{LN} (\Psi_L^m)^T \mathbf{d}^m.$	739
678		740
679		741
680	We emphasize that the above regression is ready to be computed in parallel: we can compute simultaneously the matrices	742
681	$A_L^m$ and $b_L^m$ for different trajectories. The size of the matrices $A_L^m$ is $n \times n$ , and there is no need to read and store all the data	743
682	at once, thereby dramatically reducing memory usage.	744

745 **B. Well-conditioning from coercivity.** We show next that the coercivity condition implies that  $A_{L,M}$  is well-conditioned and  
 746 positive definite for large  $M$ . More specifically, the coercivity constant provides a lower bound on the smallest singular value of  
 747  $A_{L,M}$ , provided the basis for the hypothesis space is well-conditioned (e.g. orthonormal), therefore enabling control of the  
 748 condition number of the regularized problem.  
 749 Recall the bilinear functional  $\langle\langle \cdot, \cdot \rangle\rangle$  defined in Eq. (8).  
 750 **Proposition 2.1.** Assume that the coercivity condition holds on  $\mathcal{H}_n \subset L^\infty([0, R])$  with  $c_{L,N,\mathcal{H}} > 0$ . Let  $\{\psi_1, \dots, \psi_n\}$  be a  
 751 basis of  $\mathcal{H}_n$  such that  
 752 
$$\langle \psi_p(\cdot), \psi_{p'}(\cdot) \rangle_{L^2(\rho_T^L)} = \delta_{p,p'}, \|\psi_p\|_\infty \leq S_0$$
 [17]  
 753 and  $A_{L,\infty} = (\langle\langle \psi_p, \psi_{p'} \rangle\rangle)_{p,p'} \in \mathbb{R}^{n \times n}$ . Then the smallest singular value of  $A_{L,\infty}$  satisfies  
 754 
$$\sigma_{\min}(A_{L,\infty}) \geq c_{L,N,\mathcal{H}}.$$
  
 755 Moreover,  $A_{L,\infty}$  is the a.s. limit of  $A_{L,M}$  in Eq. (16). Therefore, for large  $M$ , the smallest singular value of  $A_{L,M}$  satisfies  
 756 
$$\sigma_{\min}(A_{L,M}) \geq 0.9c_{L,N,\mathcal{H}}$$
  
 757 with probability at least  $1 - 2n \exp(-\frac{c_{L,N,\mathcal{H}}^2 M}{200n^2 c_1^2 + \frac{10c_{L,N,\mathcal{H}} c_1}{3} n})$ , where  $c_1 = R^2 S_0^2 + 1$ .  
 758 *Proof.* For each  $\mathbf{a} \in \mathbb{R}^n$ ,  
 759 
$$\mathbf{a}^T A_{L,\infty} \mathbf{a} = \langle\langle \sum_{p=1}^n a_p \psi_p, \sum_{p=1}^n a_p \psi_p \rangle\rangle \geq c_{L,N,\mathcal{H}} \left\| \sum_{p=1}^n a_p \psi_p(\cdot) \right\|_{L^2(\rho_T^L)}^2 = c_{L,N,\mathcal{H}} \|\mathbf{a}\|^2.$$
  
 760 This proves the desired bound on the smallest singular value.  
 761 Going back to the case of finite  $M$ : by the law of large numbers, the matrix  $A_{L,M} = \sum_{m=1}^M A_L^m$  converges to  $A_{L,\infty} = \mathbb{E}[A_L^m]$   
 762 as  $M \rightarrow \infty$ . Hence if the sample size  $M$  is large enough, then we apply the matrix Bernstein inequality to get the probability  
 763 estimates for the event that  $\sigma_{\min}(A_{L,M})$  is bounded below by  $0.9c_{L,N,\mathcal{H}}$ .  $\square$   
 764 **Remark 2.2.** Proposition 2.1 highlights the importance of choosing basis functions to be linearly independent in  $L^2(\rho_T^L)$  instead  
 765 of in  $L^\infty([0, R])$  for the hypothesis space  $\mathcal{H}_n$  (orthonormality can be easily obtained through Gram-Schmidt orthogonalization  
 766 if the functions are linearly independent). To see this, consider a set of basis functions consisting of piecewise polynomials  
 767 that are supported on a partition of the interval  $[0, R]$ . These functions are linearly independent in  $L^\infty([0, R])$ , but can be  
 768 linearly dependent in  $L^2(\rho_T^L)$  if some of the partitioned intervals have zero probability under the measure  $\rho_T^L$ . This would lead  
 769 to an ill-conditioned normal matrix  $A_{L,\infty}$ . This issue can deteriorate in practice when the unknown  $\rho_T^L$  is replaced by the  
 770 empirical measure  $\rho_T^{L,M}$ . In this work we use piecewise polynomials on a partition of the support of  $\rho_T^{L,M}$ , which are orthogonal  
 771 in  $L^2(\rho_T^{L,M})$ .  
 772 **C. First Order Heterogeneous Agent Systems.** For these systems the empirical error to be minimized is as in (9) in the main  
 773 text:  
 774 
$$\frac{1}{LM} \sum_{l,m,i=1}^{L,M,N} \frac{1}{N_{k_i}} \left\| \dot{x}_i^m(t_l) - \sum_{i'=1}^N \frac{1}{N_{k_{i'}}} \varphi_{k_i k_{i'}}(r_{i,i'}^m(t_l)) \mathbf{r}_{i,i'}^m(t_l) \right\|^2,$$
  
 775 over all possible  $\varphi = \{\varphi_{kk'}\}_{k,k'=1}^K \in \mathcal{H}$ . Here  $\mathbf{r}_{i,i'}(t_l)$  and  $r_{i,i'}(t_l)$  are as in Eq. (14). When given observation data,  
 776  $\{\mathbf{x}_i^m(t_l)\}_{i=1, m=1, l=1}^{N,M,L}$ , but no derivative information, we approximate the derivatives using backward differencing scheme for  
 777  $1 \leq l \leq L$  (assuming observations at  $t_0$ ); in either case we assemble the derivative vector  $\mathbf{d}$  similarly to Eq. (15), but with the  
 778 normalization  
 779 
$$\mathbf{d}^m(l) = N_{k_i}^{-1/2} \Delta \mathbf{x}_i^m(t_l) \in \mathbb{R}^d.$$
  
 780 Proceeding analogously to the homogeneous agent case, we search for  $\varphi_{kk'}$  in a  $n_{kk'}$ -dimensional hypothesis space  $\mathcal{H}_{n_{kk'}}$ ,  
 781 with basis  $\{\psi_{kk',p}\}_{p=1}^{n_{kk'}}$ , and write  $\varphi_{kk'}(r) = \sum_{k,k'=1}^K \sum_{p=1}^{n_{kk'}} a_{kk',p} \psi_{kk',p}(r)$  for some vector of coefficients  $(a_{kk',p})_{p=1}^{n_{kk'}}$ . For the  
 782 learning matrix  $\Psi_L^m$ , we will divide the columns into  $K^2$  regions, each region indexed by the pair  $(k, k')$ , with  $k, k' = 1, \dots, K$ .  
 783 We adopt the usual lexicographic partial ordering on these pairs. The columns of  $\Psi_L^m$  corresponding to  $(k, k')$  are given by  
 784 
$$\Psi_L^m(l, \tilde{n}_{kk'} + p) = N_{k_i}^{-1/2} \sum_{i' \in C_{k'}} \frac{1}{N_{k'}} \psi_{kk',p}(r_{i,i'}^m(t_l)) \mathbf{r}_{i,i'}^m(t_l) \in \mathbb{R}^d,$$
  
 785 for  $i \in C_k$  and  $1 \leq l \leq L$ , and  $\tilde{n}_{kk'} = \sum_{(k_1, k'_1) < (k, k')} n_{k_1 k'_1}$ . We define  
 786 
$$\mathbf{a} = (a_{11,1}, \dots, a_{11,n_{11}}; \dots; a_{KK,1}, \dots, a_{KK,n_{KK}}) \in \mathbb{R}^{d_0}$$
  
 787 with  $d_0 = \sum_{k,k'=1}^K n_{kk'}$ , to arrive at Eq. (16)

869 **D. Second Order Heterogeneous Agent Systems.** The learning problems of inferring the interactions of the  $\dot{x}_i$ 's and  $\xi_i$ 's can 931  
 870 be de-coupled. We start with the inference of the interactions on  $\dot{x}_i$ 's. Let the observations of the second order heterogeneous 932  
 871 agent system be  $\{\mathbf{x}_i^m(t_l), \dot{x}_i^m(t_l), \xi_i^m(t_l)\}_{l,i,m=1}^{L,N,M}$ . Let  $\mathbf{v}_i^m = \dot{x}_i^m$ . As usual, if velocities and/or accelerations are not observed, 933  
 872 they are approximated by a finite-difference (in time) scheme, for example 934  
 873  
 874 
$$\Delta \mathbf{v}_i^m(t_l) = \frac{\mathbf{v}_i^m(t_l) - \mathbf{v}_i^m(t_{l-1})}{t_l - t_{l-1}}, \quad \Delta \xi_i^m(t_l) = \frac{\xi_i^m(t_l) - \xi_i^m(t_{l-1})}{t_l - t_{l-1}},$$
 935  
 875  
 876  
 877 for  $1 \leq l \leq L$  and  $1 \leq i \leq N$  (assuming observations also at  $t_0$ ). For the data corresponding to the  $m^{th}$  initial condition, we 939  
 878 assemble the external influence (from interaction with the environment) vector  $\vec{F}^{m,\mathbf{v}}$  as: 940  
 879  
 880 
$$\vec{F}^{m,\mathbf{v}}(li) = N_{\xi_i}^{-1/2} F^{\mathbf{v}}(\mathbf{v}_i^m(t_l), \xi_i^m(t_l)) \in \mathbb{R}^d,$$
 942  
 881  
 882 and the approximated derivative of  $\mathbf{v}_i$ 's as 944  
 883  
 884  
 885 
$$\mathbf{d}^{m,\mathbf{v}}(li) = N_{\xi_i}^{-1/2} m_i \Delta \mathbf{v}_i^m(t_l) \in \mathbb{R}^d.$$
 946  
 886  
 887 We use a finite dimensional subspace  $\mathcal{H}_{n^E}^E$ , so that the candidate functions  $\varphi^E = \{\varphi_{kk'}^E\}_{k,k'=1}^K$  are expressed as  $\varphi^E(r) =$  949  
 888  $\sum_{k,k'=1}^K \sum_{p=1}^{n_{k,k'}^E} \alpha_{kk',p}^E \psi_{kk',p}^E(r)$ . Using the same ordering from previous discussion on the first order heterogeneous agent 950  
 889 system, we have, for a pair  $(k, k')$  learning matrix  $\Psi_{L,M}^{m,E}$  for the energy-based interaction kernel, 951  
 890  
 891  
 892 
$$\Psi_{L,M}^{m,E}(li, \tilde{n}^E + p) = N_{\xi_i}^{-1/2} \sum_{i' \in C_{k'}} \frac{1}{N_{k'}} \psi_{kk',p}^E(r_{i,i'}^m(t_l)) \mathbf{r}_{i,i'}^m(t_l),$$
 954  
 893  
 894  
 895  
 896 for  $1 \leq l \leq L$ ,  $i \in C_k$  and  $\tilde{n}^E = \sum_{(k_1, k'_1) < (k, k')} n_{k_1 k'_1}^E$ . The construction of the alignment-based learning matrix  $\Psi_{L,M}^{m,A}$  is 958  
 897 analogous: 959  
 898  
 899 
$$\Psi_{L,M}^{m,A}(li, \tilde{n}^A + p) = N_{\xi_i}^{-1/2} \sum_{i' \in C_{k'}} \frac{1}{N_{k'}} \psi_{kk',p}^A(r_{i,i'}^m(t_l)) \mathbf{r}_{i,i'}^m(t_l),$$
 961  
 900  
 901  
 902 for  $1 \leq l \leq L$ ,  $i \in C_k$  and  $\tilde{n}^A = \sum_{(k_1, k'_1) < (k, k')} n_{k_1 k'_1}^A$ . We put all the  $\alpha$ 's together into  $\mathbf{a}^E$  and  $\mathbf{a}^A$ , and further grouping them 964  
 903 into one big vector,  $\mathbf{a}^{\mathbf{v}} = \begin{pmatrix} \mathbf{a}^E \\ \mathbf{a}^A \end{pmatrix}$  and  $\Psi_{L,M}^{m,\mathbf{v}} = (\Psi_{L,M}^{m,E}, \Psi_{L,M}^{m,A})$ , we arrive at the final formulation, 965  
 904  
 905  
 906  
 907 
$$\frac{1}{LM} \sum_{m=1}^M \|\mathbf{d}^{m,\mathbf{v}} - \vec{F}^{m,\mathbf{v}} - \Psi_{L,M}^{m,\mathbf{v}} \mathbf{a}^{\mathbf{v}}\|_{\mathbb{R}^{LNd}}^2.$$
 968  
 908  
 909  
 910 As usual, we solve the associated normal equations of Eq. (16) with  $A_L^m := (\Psi_{L,M}^{m,\mathbf{v}})^{\top} \Psi_{L,M}^{m,\mathbf{v}}$  and  $b_L^m := (\Psi_{L,M}^{m,\mathbf{v}})^{\top} (\mathbf{d}^{m,\mathbf{v}} - \vec{F}^{m,\mathbf{v}})$ , 973  
 911 reducing the system size from  $(MLNd) \times (n^E + n^A)$  to  $(n^E + n^A)^2$ . 974  
 912  
 913 For the inference of the interactions on  $\xi_i$ 's, we let 975  
 914  
 915 
$$\vec{F}^{m,\xi}(li) = N_{\xi_i}^{-1/2} F^{\xi}(\xi_i^m(t_l)) \quad \text{and} \quad \mathbf{d}^{m,\xi}(li) = N_{\xi_i}^{-1/2} \Delta \xi_i^m(t_l),$$
 977  
 916  
 917 for  $1 \leq l \leq L$  and  $1 \leq i \leq N$ ; then the learning matrix  $\Psi_{L,M}^{m,\xi}$  is assembled similarly as 979  
 918  
 919  
 920 
$$\Psi_{L,M}^{m,\xi}(li, \tilde{n}^{\xi} + p) = N_{\xi_i}^{-1/2} \sum_{i' \in C_{k'}} \frac{1}{N_{k'}} \psi_{kk',p}^{\xi}(r_{i,i'}^m(t_l)) \mathbf{r}_{i,i'}^m(t_l),$$
 982  
 921  
 922  
 923 for  $1 \leq l \leq L$ ,  $i \in C_k$ , and  $\tilde{n}^{\xi} = \sum_{(k_1 k'_1) < (k, k')} n_{k_1 k'_1}^{\xi}$ . We then arrive at the Least Squares problem 985  
 924  
 925  
 926  
 927 
$$\frac{1}{LM} \sum_{m=1}^M \|\mathbf{d}^{m,\xi} - \vec{F}^{m,\xi} - \Psi_{L,M}^{m,\xi} \mathbf{a}^{\xi}\|_{\mathbb{R}^{LNd}}^2$$
 988  
 928  
 929  
 930 and solve it from the associated normal equations. 991  
 931  
 932  
 933  
 934  
 935

993 **E. The Final Algorithm.** Given observation data,  $\{\mathbf{x}_i^m(t_l)\}$  and  $\dot{\mathbf{x}}_i^m(t_l)$  and/or  $\xi_i^m(t_l)\}_{l,i,m=1}^{L,N,M}$ , we use the Algorithm 1 to find the  
 994 estimators for the interaction kernels.  
 995

---

**Algorithm 1** Learning Interaction Kernels from Observations  
 997 1: Input:  $\{\mathbf{x}_i^m(t_l)\}$  and/or  $\dot{\mathbf{x}}_i^m(t_l)$  and/or  $\xi_i^m(t_l)\}_{l,i,m=1}^{L,N,M}$ .  
 998 2: Output: estimators for the interaction kernels.  
 999 3: **if** First Order System **then**  
 1000 4:     Find out the maximum interaction radii  $R_{kk'}$ 's.  
 1001 5:     Construct the basis,  $\psi_{kk',p}$ 's.  
 1002 6:     Assemble the normal equations (16) (in parallel) and solve for  $\mathbf{a}$ .  
 1003 7:     Assemble  $\hat{\phi}(r) = \sum_{k,k'=1}^K \sum_{p=1}^{n_{kk'}} a_{kk',p} \psi_{kk',p}(r)$ .  
 1004 8: **else if** Second Order System **then**  
 1005 9:     Find out the maximum interaction radii  $R_{kk'}$ 's.  
 1006 10:     Construct the basis,  $\psi_{kk',p}^E$ 's and  $\psi_{kk',p}^A$ 's.  
 1007 11:     Assemble the normal equations (16) (in parallel), solve for  $\mathbf{a}^v$ , and partition it to  $\mathbf{a}^E$  and  $\mathbf{a}^A$ .  
 1008 12:     Assemble  $\hat{\phi}(r)^E = \sum_{k,k'=1}^K \sum_{p=1}^{n_{kk'}^E} a_{kk',p}^E \psi_{kk',p}^E(r)$  and  $\hat{\phi}(r)^A = \sum_{k,k'=1}^K \sum_{p=1}^{n_{kk'}^A} a_{kk',p}^A \psi_{kk',p}^A(r)$ .  
 1009 13:     **if** If there are  $\xi_i$ 's **then**  
 1010 14:         Construct the basis,  $\psi_{kk',p}^\xi$ 's.  
 1011 15:         Assemble the normal equations and solve for  $\mathbf{a}^\xi$ .  
 1012 16:         Assemble  $\hat{\phi}(r)^\xi = \sum_{k,k'=1}^K \sum_{p=1}^{n_{kk'}^\xi} a_{kk',p}^\xi \psi_{kk',p}^\xi(r)$ .  
 1013 17:     **end if**.  
 1014 18: **end if**.  
 1015

---

1016 **F. Computational Complexity.** The computational complexity is driven by the construction and solution of the least squares  
 1017 problem in Algorithm 1. Though the observation data  $\{\mathbf{x}_i^m(t_l), \dot{\mathbf{x}}_i^m(t_l), \xi_i^m(t_l)\}_{l,i,m=1}^{L,N,M}$  requires an array of size  $MLN(2d+1)$ ,  
 1018 the linear system to be solved, i.e. the system consisting of normal equations, is only of size  $n^E + n^A$ ; in the case of choosing  
 1019 the optimal basis,  $n^E$  and  $n^A$  behave like  $\mathcal{O}(M^{\frac{1}{2s+1}})$ . When the system of the normal equations is ill-conditioned or ill-posed,  
 1020 a truncated singular value decomposition will be used, which performs a singular value decomposition of the matrix  $A_{L,M}$ , and  
 1021 keeps those singular values which are above a (preset) threshold, then assemble an approximated matrix with the truncated  
 1022 singular value matrices.  
 1023

1024 Furthermore, since the  $M$  trajectories are independent, we can construct  $\Psi^{m,E}$  and other related quantities for each  
 1025 trajectory at a time (which can be done in a parallel environment with two communication needed, one to send/receive the  
 1026 maximum interaction radii's, and the other to send/receive  $A_L^m$  and  $b_L^m$  in the normal equations after they are built on the  
 1027 master core), each requires a total memory of  $LNd(n^E + n^A) + LNd + LNd$ , which is  $\mathcal{O}(LNd)$ , since  $n^E + n^A \ll LNd$ .  
 1028

1029 The computing time of the algorithm depends heavily on the time to assemble normal equations from  $M$  trajectories, which  
 1030 is  $\mathcal{O}((n^E + n^A)^2 LN^2)$ ; solving the final linear system requires time  $\mathcal{O}((n^E + n^A)^3) = \mathcal{O}(M^{\frac{3}{2s+1}})$  in the worst case, for example  
 1031 when using a highly stable truncated singular value decomposition solver.  
 1032

1033 Therefore, the algorithm is effective at inferring the interactions from a wide variety of systems; the results will be discussed  
 1034 in the next section.  
 1035

### 3. Examples

1036 We consider here four important examples of self-organized dynamics: the opinion dynamics, the particle system with the  
 1037 Lennard-Jones potential, the predator-swarm system and the phototaxis dynamics. We describe here in detail how the numerical  
 1038 simulations are set up for each of these examples. In all but the Lennard-Jones system, we set up the experiments using the  
 1039 parameters as shown in Table S1. We consider the regime with a rather small number of observations in terms of both  $M$  and  
 1040  $L$  to emphasize that our technique can achieve good results even when a relatively small number of samples is given.  
 1041

**Table S1**

$N$	# Trials	$M_{\rho_T^L}$	$[0, T_f]$
10	10	2000	$[0, cT]$

1042 Parameters used in all the examples but the Lennard-Jones system. Here the observation time  $T$  is system-specific.  $c = 2$  in all examples unless  
 1043 otherwise specified.  
 1044

1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052 We use a large number  $M_{\rho_T^L}$  (in particular,  $M_{\rho_T^L} \gg M$ ) of independent trajectories (not to be used elsewhere) to obtain an  
 1053 accurate approximation of the unknown probability measure  $\rho_T^L$  in (4) in the main text. In what follows, to keep the notation  
 1054 from becoming cumbersome, we denote by  $\rho_T^L$  this empirical approximation to  $\rho_T^L$ . We run the dynamics over the time  $[0, T_f]$   
 1055

with  $M$  different initial conditions (drawn from the dynamics-specific probability measure  $\mu_0$ ), and the observations consist of the state vector, with no derivative information, at  $L$  equidistant time samples in the time interval  $[0, T]$ . We report the relative (i.e. normalized by the norm of the true interaction kernel) error of our estimators in the  $L^2(\rho_T^L)$  norm. In the spirit of Proposition (3.4) in the main text, we also report on the error on trajectories  $\mathbf{X}(t)$  and  $\hat{\mathbf{X}}(t)$  generated by the system with the true interaction kernel and with the learned interaction kernel, on both the “training” time interval  $[0, T]$  and on a “prediction” time interval  $[T, T_f]$  ( $T_f = 2T$  unless otherwise specified), with both the same initial conditions as those used for training, and on new initial conditions (sampled according to the specified measure  $\mu_0$ ). The trajectory error will be estimated using  $M$  trajectories (we report mean and standard deviation of the error). We run a total of 10 independent learning trials and compute the mean and standard deviation of the corresponding estimators, their errors, and the trajectory errors just discussed. Since each learning trial generates different mean and standard deviation of the trajectory errors over different Initial Conditions (ICs), we also report the mean and standard deviation over the 10 learning trials for  $\text{mean}_{IC}$  and  $\text{std}_{IC}$ .

All ODE systems are evolved using **ode15s** in MATLAB® with a relative tolerance at  $10^{-5}$  and absolute tolerance at  $10^{-6}$ . We choose the finite-dimensional hypothesis space  $\mathcal{H}_n$  (with  $n$  chosen differently in each example, based on sample size) as the span of either piecewise constant or piecewise linear functions on  $n$  intervals forming a uniform partition of  $[0, R_{k,k'}]$ , where  $R_{k,k'}$  is the maximum observed pairwise distance between agents of type  $k'$  and agents in type  $k$  for  $t \in [0, T]$ .

Learning results are showcased in Fig. 5 in the main text . The first one compares the learned interaction kernel(s) to the true interaction kernel(s) (with mean and standard deviation over the total number of learning trials) with the background showing the comparison of  $\rho_T^L$  (computed on  $M_{\rho_T^L}$  trajectories, as described above) and  $\rho_T^{L,M}$  (generated from the observed data consisting of  $M$  trajectories). The second plot compares the true trajectories (evolved using the true interaction law(s)) and learned trajectories (evolved using the learned interaction law(s)) over two different sets of initial conditions – one taken from the training data, and one new, randomly generated from  $\mu_0$ . The third plot compares the true trajectories and the trajectories generated with the estimated interaction kernel, but for a different system with number of agents  $N_{\text{new}} = 4N$ , again over two different sets of randomly chosen initial conditions. Measurements of performance are also shown alongside the figures: ( $L^2(\rho_T^L)$  errors, trajectory errors, etc. Let  $\mathbf{X}(t)$  and  $\hat{\mathbf{X}}(t)$  be two sets of continuous-time trajectories; the max-in-time error is defined as

$$\|\mathbf{X} - \hat{\mathbf{X}}\|_{TM([0,T])} = \sup_{t \in [0,T]} \|\mathbf{X}(t) - \hat{\mathbf{X}}(t)\|_{\mathcal{S}}. \quad [18]$$

For second order systems with the auxiliary environment variable  $\xi_i$ 's, we are also interested in the trajectories of  $\xi_i$ , for which we may use  $\|\Xi - \hat{\Xi}\|_{TM([0,T])} = \sup_{t \in [0,T]} \|\Xi(t) - \hat{\Xi}(t)\|_{\mathcal{S}}$ .

Finally, for each example we consider adding noise to the observations: in the case of additive noise the observations are  $\{(\mathbf{X}^m(t_l) + \eta_{1,l,m}, \dot{\mathbf{X}}^m(t_l)) + \eta_{2,l,m}\}_{l=1,m=1}^{L,M}$ , while in the case of multiplicative noise they are  $\{(\mathbf{X}^m(t_l) \cdot (1 + \eta_{1,l,m}), \mathbf{X}^m(t_l) \cdot (1 + \eta_{2,l,m}))\}_{l=1,m=1}^{L,M}$ , where in both cases  $\eta_{1,l,m}$  and  $\eta_{2,l,m}$  are i.i.d. samples from a distribution modeling noise, which we will pick to be  $\text{Unif}([- \sigma, \sigma])$ . Note that in both these cases velocities are part of our observations, since with noise added in the position the inference of velocities becomes problematic due to the amplification of the noise that a simple finite difference scheme would incur.

Finally, for several examples we also report the behavior of the relative error of the estimator as a function of the number of samples  $L$  in time and of the number of trajectories  $M$ . We observe the decrease in error as  $L$  increases, which is expected but is not captured by the estimate in Thm. (3.3) in the main text. These plots are qualitatively the same for all the experiments.

We devote the next sections to the various examples, discussing setups particular to each example and corresponding results.

**A. Opinion Dynamics.** Modeling using self-organized dynamics has seen successful applications in studying and analyzing how the opinions of people influence each other and how consensus is formed based on different kinds of influence functions. We refer to these systems as opinion dynamics. We consider the first order model in Eq. (1), and the interaction kernel defined as

$$\phi(r) = \begin{cases} 1, & 0 \leq r < \frac{1}{\sqrt{2}}, \\ 0.1, & \frac{1}{\sqrt{2}} \leq r < 1, \\ 0, & 1 \leq r. \end{cases}$$

In this context  $\phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  is sometimes referred to as the scaled influence function, modeling the change of each agents' opinion by relative differences in the opinions of the other agents. Here  $\mathbf{x}_i \in \mathbb{R}^d$  is the vector opinions of agent  $i$ . Here  $\|\cdot\|$  can be taken as the normal Euclidean norm, but other metrics depending on the problem at hand may be used as well, with no changes in our definitions and constructions. The time-discretization of this system is referred to as the classical Krause model for opinion dynamics. With the specific  $\phi$  above, there is only attraction present in the system, the opinions of the agents merge into clusters, with the number of clusters significantly smaller than the number of agents. This clustering behavior severely reduces the amount of effective samples of pairwise distance observable at large times. We consider the system and test parameters given in Table S2.

**Table S2**

$d$	$M$	$L$	$T$	$\mu_0$	$n$	$\deg(\psi)$
1	50	200	10	$\mathcal{U}([0, 10]^2)$	200	0

(OD) Parameters for the system

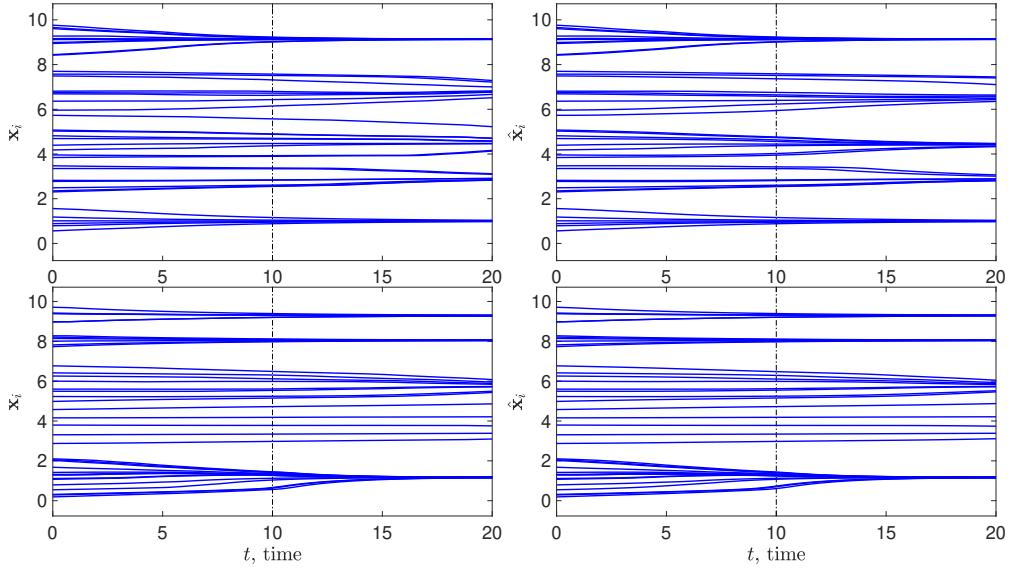


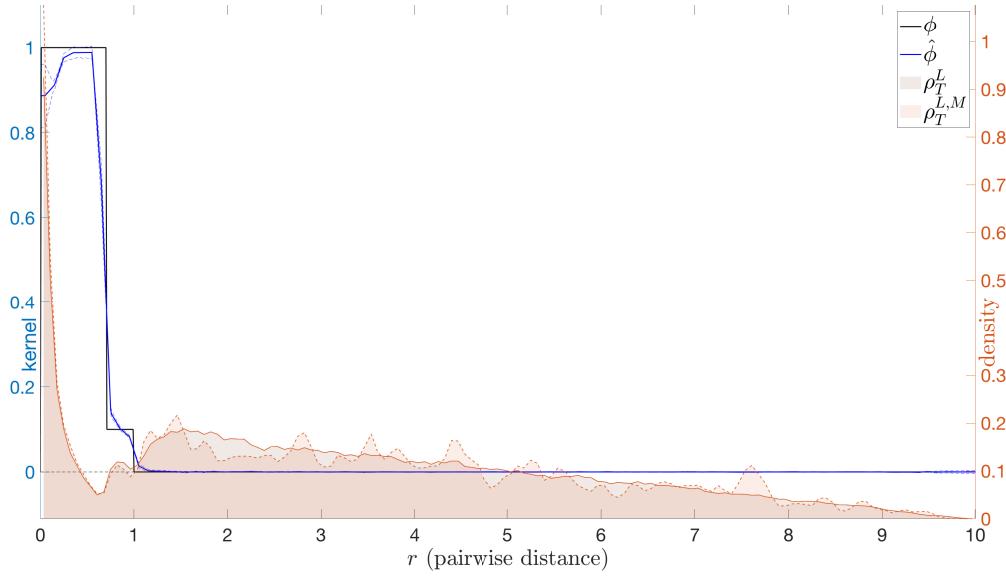
Fig. S1. (OD) Trajectories  $\mathbf{X}(t)$  and  $\hat{\mathbf{X}}(t)$  obtained with  $\phi$  and  $\hat{\phi}$  respectively, for dynamics with larger  $N_{\text{new}} = 4N$ , over two different sets of initial conditions. We are able to accurately predict the clusters (number and location). Errors are reported in Table S3.

**Table S3**

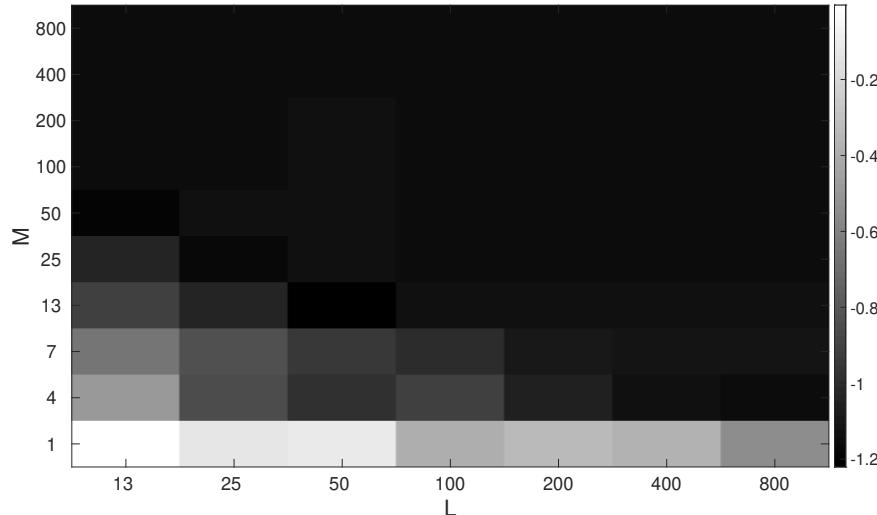
	$[0, T]$	$[T, T_f]$
$\text{mean}_{\text{IC}}$ : Training ICs	$3.5 \cdot 10^{-2} \pm 8.1 \cdot 10^{-3}$	$4.8 \cdot 10^{-2} \pm 1.4 \cdot 10^{-2}$
$\text{std}_{\text{IC}}$ : Training ICs	$5.2 \cdot 10^{-2} \pm 1.3 \cdot 10^{-2}$	$7.6 \cdot 10^{-2} \pm 2.7 \cdot 10^{-2}$
$\text{mean}_{\text{IC}}$ : Random ICs	$3.2 \cdot 10^{-2} \pm 7.4 \cdot 10^{-3}$	$4.6 \cdot 10^{-2} \pm 1.2 \cdot 10^{-2}$
$\text{std}_{\text{IC}}$ : Random ICs	$5.0 \cdot 10^{-2} \pm 1.7 \cdot 10^{-2}$	$7.2 \cdot 10^{-2} \pm 2.7 \cdot 10^{-2}$
$\text{mean}_{\text{IC}}$ : Larger $N$	$3.1 \cdot 10^{-2} \pm 2.0 \cdot 10^{-3}$	$7.3 \cdot 10^{-2} \pm 4.1 \cdot 10^{-3}$
$\text{std}_{\text{IC}}$ : Larger $N$	$2.1 \cdot 10^{-2} \pm 2.1 \cdot 10^{-3}$	$6.1 \cdot 10^{-2} \pm 4.2 \cdot 10^{-3}$

(OD) Trajectory Errors: ICs used in the training set (first two rows), new IC's randomly drawn from  $\mu_0$  (second set of two rows), for ICs randomly drawn for a system with  $4N$  agents (last two rows). Means and std's are over 10 learning runs.

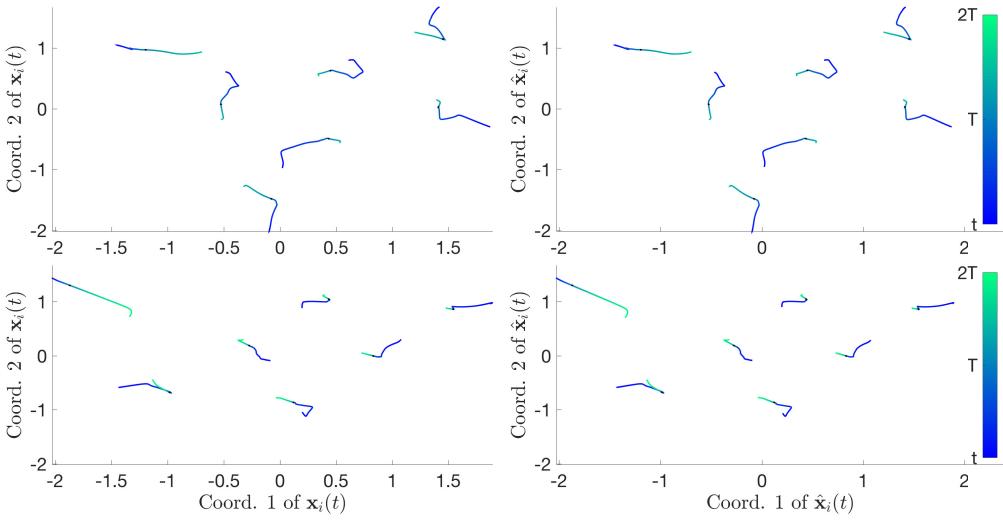
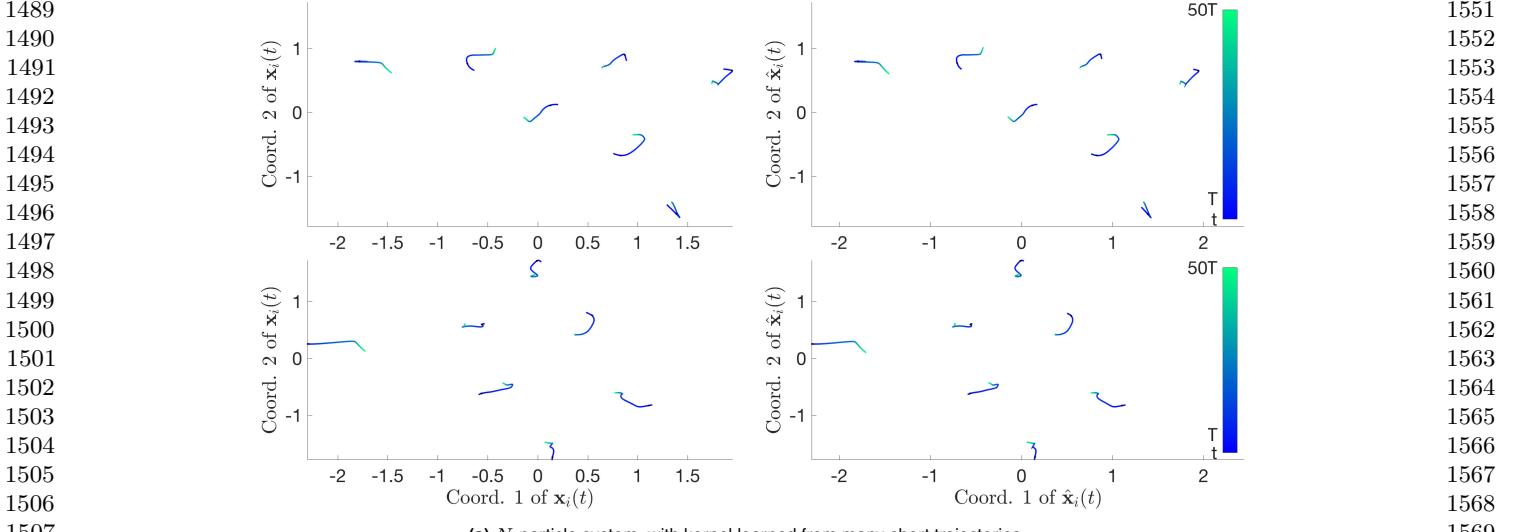
Fig. S1 shows the comparison between the estimated interaction kernel  $\hat{\phi}$  (as the mean over learning trials) and the true one,  $\phi$ . We obtain a faithful approximation of the true interaction kernel, including near the discontinuity and the compact support. Our estimator also performs well near 0, notwithstanding that information of  $\phi(0)$  is lost due to the structure of the equations, that have terms of the form  $\phi(0)\vec{0} = \vec{0}$ . The same figure also compares the trajectories generated by the system governed by  $\phi$  and that governed by  $\hat{\phi}$ . Table S3 reports the max-in-time error for those trajectories. We also test the robustness to noise, by adding noise to the observations of both positions and velocities, as described above: the estimated kernel is shown in Figure S2. Figure S3 shows the behavior of the error of the estimator as both  $L$  and  $M$  are increased.



**Fig. S2.** (OD) Interaction kernel learned with  $\text{Unif.}([-\sigma, \sigma])$  additive noise, for  $\sigma = 0.1$  in the observed positions and velocities. The estimated kernels are minimally affected, mostly in regions with small  $\rho_T^L$  and near 0.



**Fig. S3.** (OD) Relative error, in  $\log_{10}$  scale, of  $\hat{\phi}$  as a function of  $L$  and  $M$ . The error decreases both in  $L$  and  $M$ , in fact roughly in the product  $ML$ , at least when  $M$  and  $L$  are not too small.  $M = 1$  does not seem to suffice, no matter how large  $L$  is, due to the limited amount of "information" contained in a single trajectory.



**Fig. S4.** (LJ) (a) and (b) presents trajectories  $\mathbf{X}(t)$  (left) and  $\widehat{\mathbf{X}}(t)$  (right) obtained with  $\phi$  and  $\widehat{\phi}$  respectively, for initial conditions in the training dataset (top) and randomly sampled initial conditions (bottom). The time  $T$  is as in Table S5. Trajectory errors for all cases are reported in Table S7.

**B. Interacting Particles in Lennard-Jones Potential.** The expression of the Lennard-Jones potential is

$$\Phi(r) = 4\epsilon \left[ \left(\frac{\sigma}{r}\right)^{12} - \left(\frac{\sigma}{r}\right)^6 \right] = \epsilon \left[ \left(\frac{r_m}{r}\right)^{12} - 2 \left(\frac{r_m}{r}\right)^6 \right]$$

where  $\epsilon$  is the depth of the potential well,  $\sigma$  is the finite distance at which the inter-particle potential is zero,  $r$  is the distance between the particles, and  $r_m$  is the distance at which the potential reaches its minimum. At  $r_m$ , the potential function has the value  $-\epsilon$ . The  $r^{-12}$  term describes Pauli repulsion at short ranges due to overlapping electron orbitals, and the  $r^{-6}$  term describes attraction at long ranges (van der Waals force, or dispersion force). We set  $\epsilon = 10$  and  $\sigma = 1$  in our simulations.

In the experiments, whose results are represented in Fig. 1 in the main text, the distribution  $\mu_0$  for the  $M$  i.i.d. initial conditions is a standard Gaussian vector in  $\mathbb{R}^{2N}$ . In this Lennard-Jones interacting system, one has to be careful in choosing the observation time interval. Since the minimum distance between the particles at initial configurations is very close to 0 with high probability, the particles have very large velocities (e.g.  $\sim 10^{22}$ ) due to the singularity of the interaction kernel at 0. This obstruction made the learning algorithm infeasible since our algorithm is for learning bounded kernels. Therefore, we chose an observation time starting from a suitable time  $t_0$ , small but positive. On the other side of the training time interval, since the system evolves to equilibrium configurations very quickly, we observe the dynamics up to a time  $T$  which is a fraction of the equilibrium time. In each sampling regime, we observe the dynamics at discrete times  $\{t_i\}_{i=2,\dots,L}$  and then use the standard finite difference method to obtain a faithful approximation of velocities of agents.

1613  
1614  
1615  
1616  
1617  
1618  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673  
1674  
**Table S4**  

$N$	$d$	$\mu_0$	# Trials	$M_{\rho_T^L}$	$[t_0, T_f]$	$\deg(\psi_{kk'})$
7	2	$N(0, I_{2N})$	10	2000	$[t_0, cT]$	1

(LJ) Parameters used in Lennard-Jones system

  
**Table S5**  

	$M$	$L$	$n$	$[t_0, T]$	$c$
Many short traj.	200	91	600	$[0.001, 0.01]$	50
Single long traj.	20	4991	600	$[0.001, 0.5]$	2

(LJ) Observation parameters for the Lennard-Jones system

  
**Table S6**  

	Many short trajectories	a few long trajectories
Rel. Err. for $\hat{\phi}$	$6.6 \cdot 10^{-2} \pm 5 \cdot 10^{-3}$	$7.2 \cdot 10^{-2} \pm 1 \cdot 10^{-2}$

(LJ) Relative error of the estimator for the Lennard-Jones system

The estimator belongs to a piecewise linear function space  $\mathcal{H}_n$  of dimension  $n = 600$ . As reported in Fig.1 of the main text, the estimated interaction kernel  $\hat{\phi}$  approximates the true interaction kernel  $\phi$  well in the regions where  $\rho_T^L$  (and  $\rho_T$ ) is large, i.e. regions with an abundance of observed values of pairwise distances to reconstruct the interaction kernel. The dependency on  $T$  of  $\rho_T^L$ , and of the space  $L^2(\rho_T^L)$  (see (5) in the main text) used for learning, is rather pronounced, as may be seen from the histogram visualization also in Fig. 1. As usual we also compare trajectories  $\hat{\mathbf{X}}(t)$  generated by the system with the estimated interaction kernel learned with trajectories  $\mathbf{X}(t)$  generated by the original system, given the same initial conditions at  $t_0$ , both on the learning interval  $[t_0, T]$  and on larger time intervals  $[t_0, cT]$ . Figure S4 provides a visualization of such trajectories. Visualization of the corresponding systems with a larger number of agents  $N_{\text{new}}$  can be found in Figure 1 of the main text. We report the estimation errors of the interaction kernel and the trajectory errors in Tables S6 and S7.

Table S6 shows the mean and standard deviations of the relative  $L^2(\rho_T)$  errors of the kernel estimators in 10 different simulations. We report the relative errors of trajectory prediction in SI Sec.3B.

  
**Table S7**  

	$[t_0, T]$	$[T, T_f]$
mean <sub>IC</sub> : Training ICs	$1.6 \cdot 10^{-3} \pm 2 \cdot 10^{-4}$	$1.7 \cdot 10^{-2} \pm 2 \cdot 10^{-3}$
std <sub>IC</sub> : Training ICs	$4.6 \cdot 10^{-4} \pm 5 \cdot 10^{-5}$	$2.1 \cdot 10^{-2} \pm 4 \cdot 10^{-3}$
mean <sub>IC</sub> : Random ICs	$1.6 \cdot 10^{-3} \pm 2 \cdot 10^{-4}$	$1.7 \cdot 10^{-2} \pm 2 \cdot 10^{-3}$
std <sub>IC</sub> : Random ICs	$4.5 \cdot 10^{-4} \pm 5 \cdot 10^{-5}$	$1.9 \cdot 10^{-2} \pm 2 \cdot 10^{-3}$
mean <sub>IC</sub> : Larger $N$	$6.2 \cdot 10^{-2} \pm 7 \cdot 10^{-3}$	$6.2 \cdot 10^{-2} \pm 2 \cdot 10^{-2}$
std <sub>IC</sub> : Larger $N$	$8.2 \cdot 10^{-3} \pm 7 \cdot 10^{-4}$	$3.0 \cdot 10^{-2} \pm 1 \cdot 10^{-2}$
mean <sub>IC</sub> : Training ICs	$3.4 \cdot 10^{-3} \pm 1 \cdot 10^{-3}$	$5.1 \cdot 10^{-3} \pm 2 \cdot 10^{-3}$
std <sub>IC</sub> : Training ICs	$2.7 \cdot 10^{-3} \pm 2 \cdot 10^{-3}$	$6.6 \cdot 10^{-3} \pm 3 \cdot 10^{-3}$
mean <sub>IC</sub> : Random ICs	$4.1 \cdot 10^{-3} \pm 2 \cdot 10^{-3}$	$8.7 \cdot 10^{-3} \pm 8 \cdot 10^{-3}$
std <sub>IC</sub> : Random ICs	$3.6 \cdot 10^{-3} \pm 2 \cdot 10^{-3}$	$1.5 \cdot 10^{-2} \pm 2 \cdot 10^{-2}$
mean <sub>IC</sub> : Larger $N$	$7.7 \cdot 10^{-2} \pm 1 \cdot 10^{-2}$	$6.6 \cdot 10^{-2} \pm 3 \cdot 10^{-2}$
std <sub>IC</sub> : Larger $N$	$1.5 \cdot 10^{-2} \pm 1 \cdot 10^{-2}$	$5.7 \cdot 10^{-2} \pm 3 \cdot 10^{-2}$

(LJ) Trajectory Errors for Many Short Trajectories Learning (top) and Single Large Time Trajectories Learning (bottom)

We also test the convergence of our estimator as  $M \rightarrow \infty$ : we choose the parameters for observations and learning as in Table S8. It is important that we choose the dimension  $n$  of hypothesis space to be dependent on  $M$ , as dictated by Thm. (3.3) in the main text. Also, in this experiment (and this experiment only!) we observe the true derivatives (instead of approximating them by finite differences of positions), as those would introduce a bias term that does not vanishes unless  $L$  also increased with  $n$ .

14 of 27

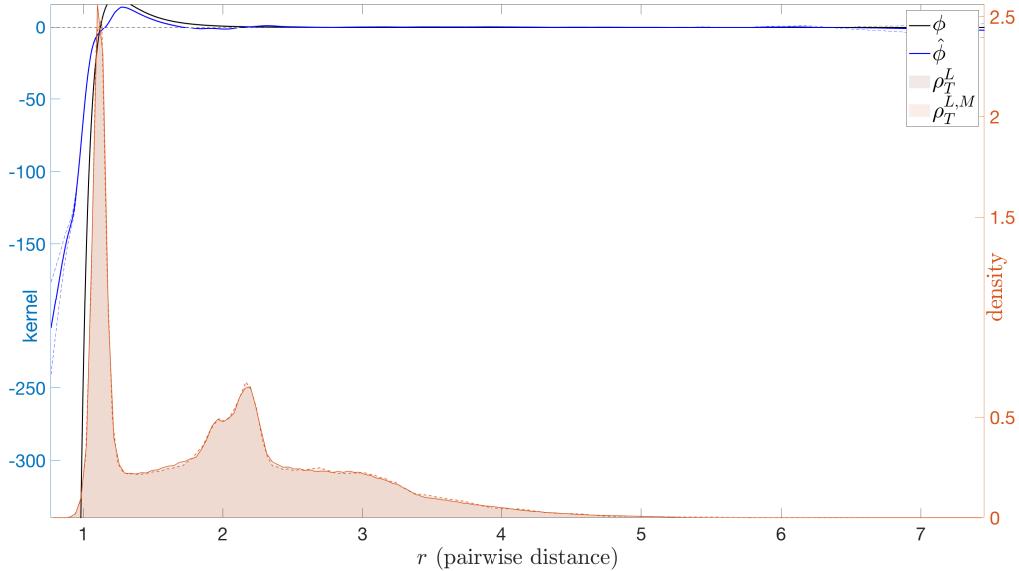
Fei Lu, Ming Zhong, Sui Tang, Mauro Maggioni

Table S8

$[t_0, T]$	$L$	$\log_2(M)$	$n$
$[0.001, 0.01]$	10	12 : 21	$64(M/\log M)^{0.2}$

(LJ) Observation parameters in the plot of convergence rate

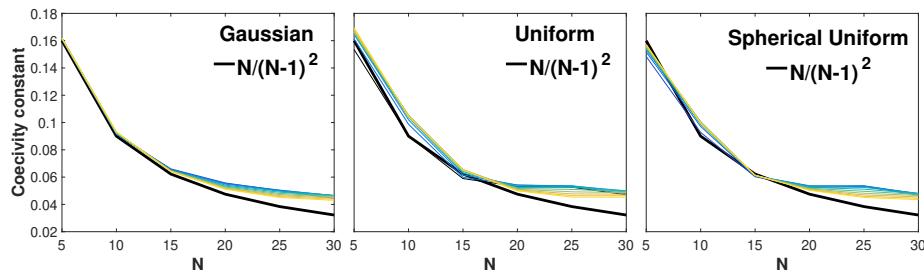
We obtain a decay rate for  $\|\hat{\phi}(\cdot) - \phi(\cdot)\|_{L^2(\rho_T^L)}$  around  $M^{-0.36}$ , which is close to the theoretical optimal learning rate  $M^{-0.4}$  – see Fig. 2 in the main text. We impute this (small) difference to the singularity of the Lennard-Jones interaction kernel at 0, which makes this interaction kernel not admissible in our learning theory.



**Fig. S5.** (LJ) Interaction kernel learned with  $\text{Unif}([-\sigma, \sigma])$  additive noise, for  $\sigma = 0.1$ , in the observed positions and observed velocities; here  $M = 500$ ,  $L = 2000$ , with all the other parameters as in Table S5.

However, the singularity of the Lennard-Jones interaction kernel at 0 forces the particles close to each other to be repel each other. Also, the system evolves rapidly to a steady-state, and the particles only explore a bounded region due to the large range attraction. Therefore, to obtain a well-supported non-degenerate measure  $\rho_T^L$ , we should make observations on a time interval that avoids reaching either the singularity of the interaction kernel or the steady-state. The restriction of the Lennard-Jones interaction kernel to the support of  $\rho_T^L$  is bounded and smooth, and hence our learning theory applies and we achieve an almost optimal rate of learning in the numerical experiments. The estimated interaction kernel with noisy observation is visualized in Figure S5.

Finally, Fig. S6 reports numerical validations of the coercivity condition in Definition 1.1 for this system. We consider the number of agents  $N$  ranging from 5 to 30, three different initial distributions  $\mu_0$ , and observations on different time intervals. The coercivity constants computed by Monte Carlo sampling are close to the theoretical lower bound in all these cases.



**Fig. S6.** (LJ) Coercivity condition validation in 2D Lennard-Jones system with different  $N$ . We compute the empirical coercivity constant  $c_{L,N,\mathcal{H}}$  defined in Eq. (7), with  $\mathcal{H}$  consisting of 200 piecewise constant basis functions with random coefficients, using  $M = 131,072$  trajectories with initial conditions drawn from  $\mu_0$ . Three initial distributions for  $\mu_0$  are tested: the standard Gaussian vector in  $\mathbb{R}^{2N}$  (left), the uniform distribution on  $[-0.5, 0.5]^{2N}$  (middle), and the uniform distributions on the unit spheres in  $\mathbb{R}^{2N}$  (right). Ten different lengths of trajectories are considered (represented in each figure by the colored curves above the black curve, the theoretical lower bound of  $c_{L,N,\mathcal{H}}$ ): each with the same initial time  $t_1 = 0.001$ , but the end time  $t_L$  ranges from 0.0059 to 0.0509 with a uniform time gap  $10^{-4}$ . In all these ten sampling regimes (all are short time periods), the coercivity constant is around  $\frac{N-1}{N^2}$ , matching the theoretical lower bound in Thm. 3.1 for one time step. We also note that  $c_{L,N,\mathcal{H}}$  appears to not go to 0 as  $N$  increases, consistent with the conjecture that in rather great generality  $c_{L,N,\mathcal{H}}$  stays bounded away from 0 independently of  $N$ .

1861 **C. Predator-Swarm system.** There is an increasing amount of literature in discussing models of self-organized animal motion 1923  
 1862 (5–15). Even more challenging is modeling interactions between agents of multiple types, in complex and emergent physical and 1924  
 1863 social phenomena (11, 16–19). We consider here a representative heterogeneous agent dynamics: a Predator-Swarm system 1925  
 1864 with a group of preys and a single predator, governed by either a first order or a second order system of ODE’s. The intensity 1926  
 1865 of interaction(s) between the single predator and group of preys can be tuned with parameters, determining dynamics with 1927  
 1866 various interesting patterns (from confusing the predator with fast preys, to chase, to catch up to one prey). Since there is 1928  
 1867 one single predator in the system, there is no predator-predator interaction to be learned. The interaction kernels (prey-prey, 1929  
 1868 predator-prey) have both short-range repulsion to prevent the agents to collide, and long-range attraction to keep the agents in 1930  
 1869 the flock. Because of the strong short-range repulsion, the pairwise distances stay bounded away from  $r = 0$ . We will see that 1931  
 1870 these difficulties, similar to those confronted with the Lennard-Jones interaction kernel, do not prevent us from learning the 1932  
 1871 interactions kernels. 1933  
 1872 1934  
 1873 In our notation for the heterogeneous system, the set  $C_1$  corresponds to the set of preys, and  $C_2$  to the set consisting of the 1935  
 1874 single predator. 1936  
 1875 1937  
 1876 **Predator-Swarm, 1<sup>st</sup> order** (PS1<sup>st</sup>). We start from the first order system. It is a special case of the first order heterogeneous 1938  
 1877 agent systems we considered, with the following interaction kernels: 1939  
 1878 1940  
 1879 1941  
 1880 1942  
 1881 1943  
 1882 1944  
 1883 1945  
 1884 1946  
 1885  $\phi_{1,1}(r) = 1 - r^{-2}, \quad \phi_{1,2}(r) = -2r^{-2}, \quad \phi_{2,1}(r) = 3r^{1.5}, \quad \phi_{2,2}(r) \equiv 0.$  1947  
 1886 1948  
 1887 1949  
 1888 1950  
 1889 1951  
 1890 1952  
 1891 1953  
 1892 The simulation parameters are given in Table S9. 1954  
 1893 1955  
 1894 1956  
 1895 1957  
 1896 1958  
 1897 1959

**Table S9**

$d$	$N_1$	$N_2$	$M$	$L$	$T$
2	9	1	50	200	5
$n_{1,1}$	$n_{1,2} = n_{2,1}$	$n_{2,2}$	$\deg(\psi_{kk'})$	Preys $\mu_0^X$	Pred. $\mu_0^X$
360	120	64	[1, 1; 1, 0]	Unif. on ring [0.5, 1.5]	Unif. on disk at 0.1

(PS1<sup>st</sup>) System parameters for first order Predator-Swarm system

1909 In the first column of Fig. 5 in the main text, we show the comparison of the learned interaction kernels versus the true 1971  
 1910 interaction kernels (with  $\rho_T^{L,kk'}$  and  $\rho_T^{L,M,kk'}$  shown in the background), and the comparison of true and learned trajectories 1972  
 1911 over two different sets of initial conditions. 1973  
 1912 1974  
 1913 1975

1914 As is shown in the top left a portion (4 sub-figures) of Fig. 5 in the main text, we are able to match faithfully all four 1976  
 1915 learned interactions to their corresponding true interactions over the range of  $\rho_T$  when the pairwise distance data is abundant. 1977  
 1916 We are not able to learn the interaction kernels for  $r$  close to 0, demonstrated by the larger area of uncertainty (surrounded by 1978  
 1917 the dashed lines) towards 0: first, the prey-to-prey interaction is preventing preys colliding into each other; second, in the case 1979  
 1918 of chasing predators, the preys are able to push away the predator. The predator-to-prey and prey-to-predator interactions are 1980  
 1919 learned over the same set of pairwise distance data, however, we are able to learn the details of the two interaction kernels, and 1981  
 1920 judging from the learned interaction kernels, they are not simply negative of each other. The predator-to-predator interaction 1982  
 1921 simply is learned as a zero function, even though there is no pairwise distance data of a predator to a different predator. Errors 1983  
 1922 in their corresponding  $L^2(\rho_T^{L,kk'})$  norms are reported in Table S10. 1984

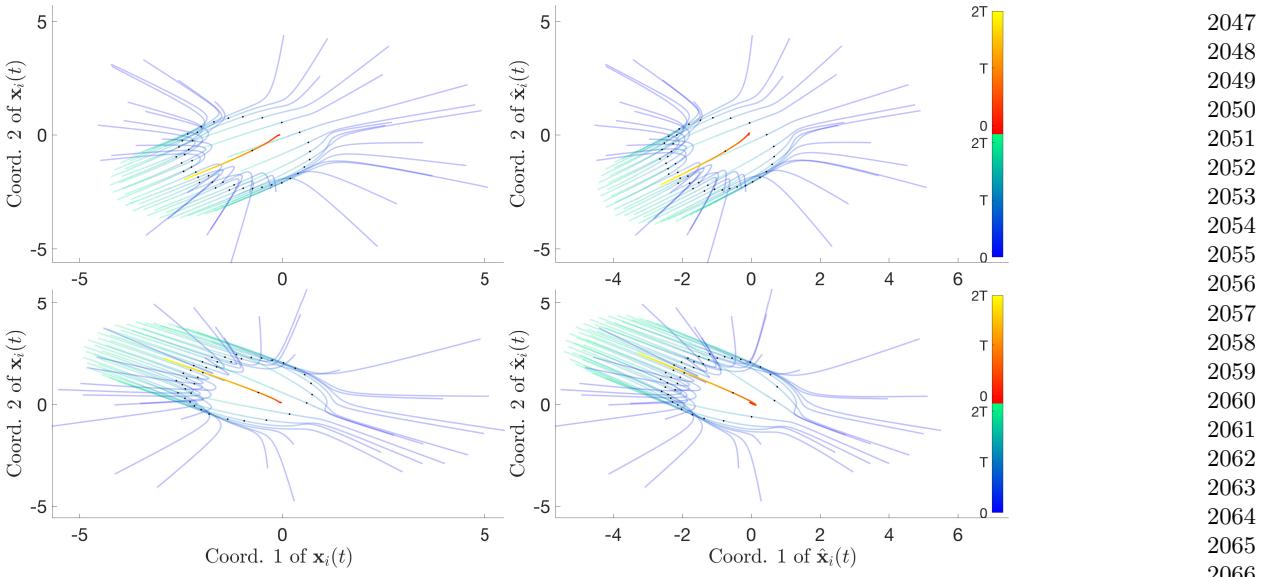


Fig. S7. (PS1<sup>st</sup>) Trajectories  $\mathbf{X}(t)$  and  $\hat{\mathbf{X}}(t)$  obtained with  $\phi$  and  $\hat{\phi}$  respectively, for two randomly chosen initial conditions and evolved for  $N_{\text{new}}$  agents (with the same setup as in the case of  $N$  agents). Trajectory errors are shown in Table S11.

The trajectory comparisons are shown in the bottom left portion (4 sub-figures) of Fig. 5 in the main text. We use color changing lines to indicate the movement of agents in time: with the blue-to-green lines attached to preys and the red-to-yellow line for the predator). The black dot on the trajectories indicate the position of the agents at time  $t = T$ , and it shows the time divide: the first half of the time,  $[0, T]$ , is used for learning; and the second half of the time,  $[T, T_f]$ , is used for prediction.

And the first row of 2 sub-figures show the comparison of the trajectories over the initial condition taken from training data, it shows (visually) no major difference between the two, except one of the prey-trajectory, is having a bigger loop in the learned trajectories. The second row of 2 sub-figures compares the trajectories from a randomly chosen initial condition (outside of the training set). We are able to predict the movement of the predator in the learned trajectories, and movement of most preys. In Fig. S7 we compare the true and predicted trajectories over a corresponding system a dynamics but with a larger number  $N_{\text{new}}$  of agents. Table S11 reports the max-in-time error Eq. (18) in the trajectories in all cases considered. We consider the effect of adding noise to observations, with results visualized in Fig. 8 of the main text.

**Table S10**

Rel. Err. for $\hat{\phi}_{1,1}$	$5.6 \cdot 10^{-2} \pm 1.1 \cdot 10^{-3}$
Rel. Err. for $\hat{\phi}_{1,2}$	$6.6 \cdot 10^{-3} \pm 2.4 \cdot 10^{-3}$
Rel. Err. for $\hat{\phi}_{2,1}$	$2.7 \cdot 10^{-2} \pm 8.9 \cdot 10^{-3}$
Abs. Err. for $\hat{\phi}_{2,2}$	0

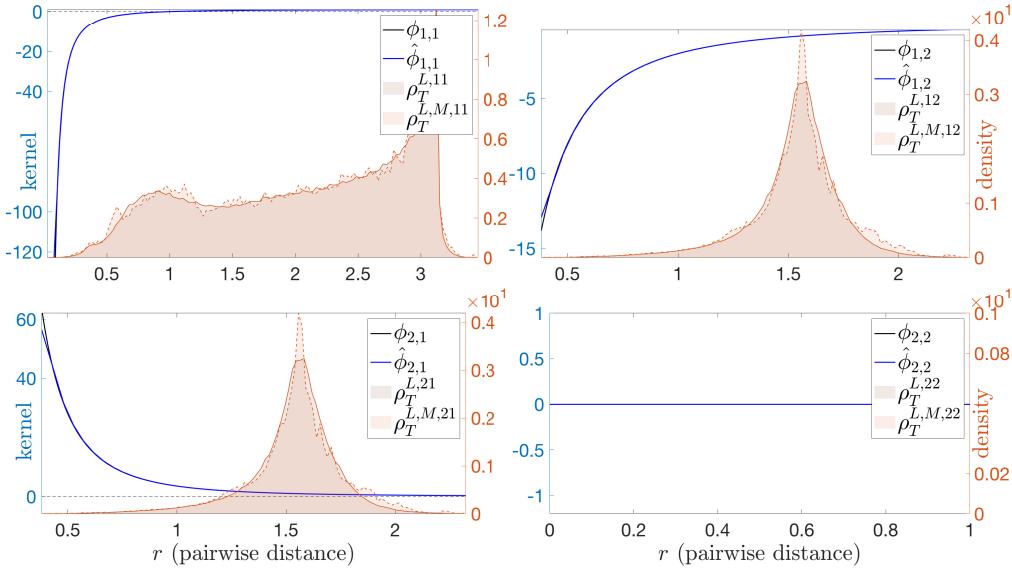
(PS1<sup>st</sup>) Estimator Errors

**Table S11**

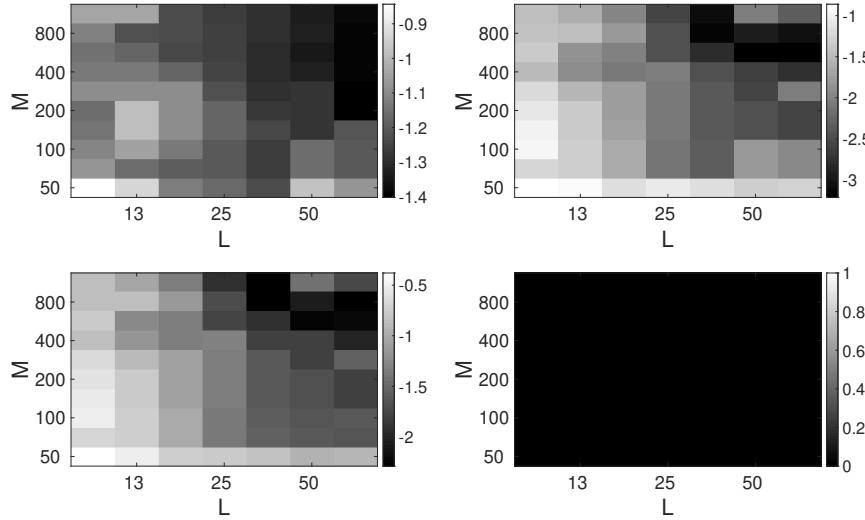
	$[0, T]$	$[T, T_f]$
mean <sub>IC</sub> : Training ICs	$4.2 \cdot 10^{-2} \pm 1.0 \cdot 10^{-2}$	$1.1 \cdot 10^{-1} \pm 3.0 \cdot 10^{-2}$
std <sub>IC</sub> : Training ICs	$7.2 \cdot 10^{-2} \pm 5.6 \cdot 10^{-2}$	$1.9 \cdot 10^{-1} \pm 1.4 \cdot 10^{-1}$
mean <sub>IC</sub> : Random ICs	$3.8 \cdot 10^{-2} \pm 1.4 \cdot 10^{-2}$	$9.5 \cdot 10^{-2} \pm 3.2 \cdot 10^{-2}$
std <sub>IC</sub> : Random ICs	$5.5 \cdot 10^{-2} \pm 6.2 \cdot 10^{-2}$	$1.4 \cdot 10^{-1} \pm 1.4 \cdot 10^{-1}$
mean <sub>L</sub> : Larger $N$	$4.2 \cdot 10^{-1} \pm 1.7 \cdot 10^{-1}$	$3.1 \pm 4.6$
std <sub>L</sub> : Larger $N$	$1.7 \cdot 10^{-1} \pm 9.6 \cdot 10^{-2}$	$15.8 \pm 27.4$

(PS1<sup>st</sup>) Trajectory Errors

We show numerically that our learning approach is robust to the choice of hypothesis space, as predicted by the theory, by testing on the Predator-Swarm, 1<sup>st</sup>-order system with the B-splines basis. Results are shown in Fig. S8. Note that the estimators perform similarly in comparison with Fig. 8 of the main text and are consistent with the error statistics in Table S11, in both of which the hypothesis space uses piece-wise polynomial basis.



**Fig. S8.** (PS1<sup>st</sup>) Comparison of interaction kernels (true versus learned) when the learned kernels are generated by linear B-splines ( $n$  as in the other case considered for this system). The relative error (in  $L^2(\rho_T)$  norm) for prey on prey interaction is:  $6.6 \cdot 10^{-2}$ ; for predatory on prey:  $6.1 \cdot 10^{-3}$ ; for prey on predator:  $3.6 \cdot 10^{-2}$ ; and finally for predator on predator: 0.



**Fig. S9.** (PS1) Relative error, in  $\log_{10}$  scale, of  $\hat{\phi}_{k,k'}^E$ , (with  $(k, k')$  increasing lexicographically from top-left to bottom-right) as a function of  $L$  and  $M$ . The error decreases both in  $L$  and  $M$ , in fact roughly in the product  $ML$ . The fourth plot is an identically 0 absolute error, because both  $\phi_{2,2}^E$  and its estimator are identically 0, since there is only one predator. Note  $M \gg 1$  seems to be needed for accurate inference of the interaction kernels, regardless of how large  $L$  is: the trajectories explored for small  $M$  do not explore enough configuration to enable estimation, suggesting that the limit  $M \rightarrow +\infty$  considered in this work is of fundamental importance, at least for non-ergodic systems.

**Predator-Swarm, 2<sup>nd</sup>-order (PS2<sup>nd</sup>).** The second order Predator-Swarm system is a special case of the second order system which is considered in this paper, without alignment-based interactions and without environment variables  $\xi_i$ 's, similar to the Cucker-Dong model of repulsion-attraction (20) and D'Orsogna-Bertozzi model for modeling fish school formation (5, 6) without the non-collective forcing term. The energy-based interactions are

$$\phi_{1,1}(r) = 1 - r^{-2}, \quad \phi_{1,2}(r) = -r^{-2}, \quad \phi_{2,1}(r) = 1.5r^{-2.5}, \quad \phi_{2,2}(r) \equiv 0.$$

The non-collective change on  $\dot{x}_i$  is  $F_i^v(\dot{x}_i, \xi_i) = -\nu_{k_i} \dot{x}_i$ , where the friction constants are type-based and  $\nu_k = 1$  for all  $k = 1, \dots, K$ ; and the mass of each agent is  $m_i = 1$  for all  $i = 1, \dots, N$ . We consider the system and test parameters given in table S12 (the initial velocity of preys and predator are fixed at  $0 \in \mathbb{R}^2$ ).

2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267  
2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294

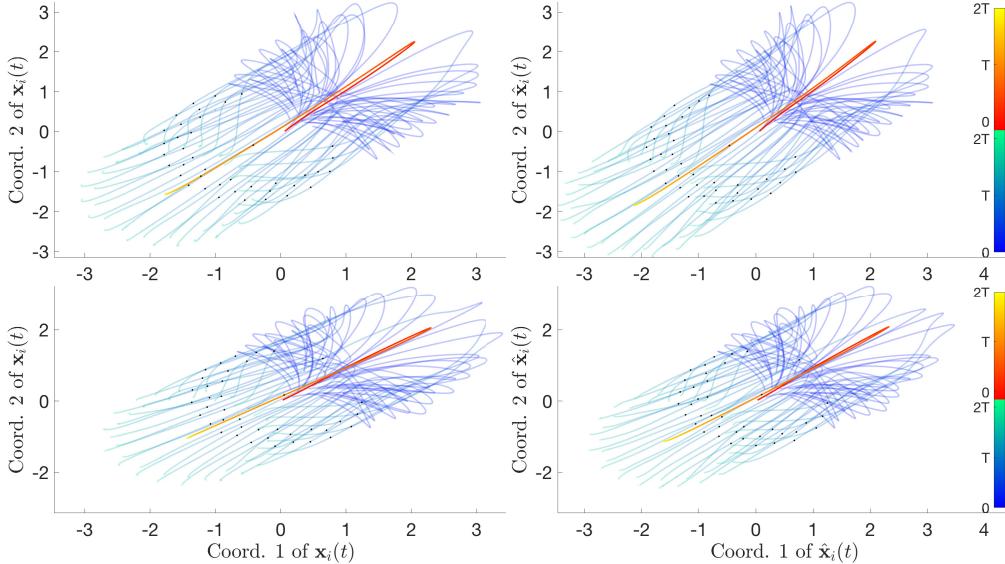
**Table S12**

$d$	$N_1$	$N_2$	$M$	$L$	$T$
2	9	1	150	300	10
$n_{1,1}$	$n_{1,2} = n_{2,1}$	$n_{2,2}$	$\deg(\psi_{kk'}^E)$	Preys $\mu_0^X$	Pred. $\mu_0^X$
1620	540	180	[1, 1; 1, 0]	Unif. on $[0.1, 1]^2$	Unif. on $[0, 0.08]^2$

(PS2<sup>nd</sup>) System Parameters

Note that the two dynamics, predator-prey 1<sup>st</sup> order and predator-prey 2<sup>nd</sup> order, use a similar set of interaction kernels, however, the resulting dynamics are significantly different from each other, as demonstrated in both the distribution of pairwise distance data and in the trajectories.

In the middle column of Fig. 5 in the main text, we show the comparison of the learned interaction kernels versus the true interaction kernels (with  $\rho_{T,r}^{L,kk'}$  and  $\rho_{T,r}^{L,M,kk'}$  shown in the background), and the comparison of true and learned trajectories over two different sets of initial conditions. Similar observations to those for the 1<sup>st</sup> order system apply here. Errors of the estimators in the  $L^2(\rho_T^{L,kk'})$  norms are reported in Table S13. The test on trajectories (bottom middle portion (4 sub-figures) of Fig. 5 in the main text) shows visually the accuracy of the predicted trajectories, quantified by the numerical report in Table S14. We also compare in Fig. S10 the true and learned trajectories over a corresponding system with  $N_{\text{new}}$  agents. We consider the effect of adding noise to observations, with results visualized in Figure S11. Figures S9 and S12 show the behavior of the error of the estimator (for systems (PS1<sup>st</sup>) and (PS2<sup>nd</sup>) respectively) as both  $L$  and  $M$  are increased.



**Fig. S10.** (PS2<sup>nd</sup>) Trajectories  $\mathbf{X}(t)$  and  $\hat{\mathbf{X}}(t)$  obtained with  $\phi$  and  $\hat{\phi}$  respectively, for two randomly chosen initial conditions and evolved for  $N_{\text{new}}$  agents (with the same setup as in the case of  $N$  agents). Trajectory errors are shown in Table S14.

2341  
2342  
2343  
2344  
2345  
2346  
2347  
2348  
2349  
2350  
2351  
2352  
2353  
2354  
2355  
2356

**Table S13**

Rel. Err. for $\hat{\phi}_{1,1}^E$	$1.5 \cdot 10^{-1} \pm 5.0 \cdot 10^{-2}$
Rel. Err. for $\hat{\phi}_{1,2}^E$	$1.3 \cdot 10^{-1} \pm 1.1 \cdot 10^{-2}$
Rel. Err. for $\hat{\phi}_{2,1}^E$	$7.1 \cdot 10^{-1} \pm 3.8 \cdot 10^{-1}$
Abs. Err. for $\hat{\phi}_{2,2}^E$	0

(PS2<sup>nd</sup>) Estimator Errors

Table S14

	$[0, T]$	$[T, T_f]$
mean <sub>IC</sub> : Training ICs	$3.5 \cdot 10^{-1} \pm 1.2 \cdot 10^{-1}$	$7.9 \cdot 10^{-1} \pm 2.1 \cdot 10^{-1}$
std <sub>IC</sub> : Training ICs	$6.5 \cdot 10^{-1} \pm 2.7 \cdot 10^{-1}$	$1.2 \pm 3.7 \cdot 10^{-1}$
mean <sub>IC</sub> : Random ICs	$3.5 \cdot 10^{-1} \pm 1.2 \cdot 10^{-1}$	$8.0 \cdot 10^{-1} \pm 2.3 \cdot 10^{-1}$
std <sub>IC</sub> : Random ICs	$5.8 \cdot 10^{-1} \pm 1.6 \cdot 10^{-1}$	$1.2 \pm 3.1 \cdot 10^{-1}$
mean <sub>IC</sub> : Larger $N$	$2.0 \cdot 10^{-1} \pm 3.0 \cdot 10^{-2}$	$4.6 \cdot 10^{-1} \pm 1.2 \cdot 10^{-1}$
std <sub>IC</sub> : Larger $N$	$1.1 \cdot 10^{-1} \pm 1.4 \cdot 10^{-2}$	$2.5 \cdot 10^{-1} \pm 5.6 \cdot 10^{-2}$

(PS2<sup>nd</sup>) Trajectory Errors

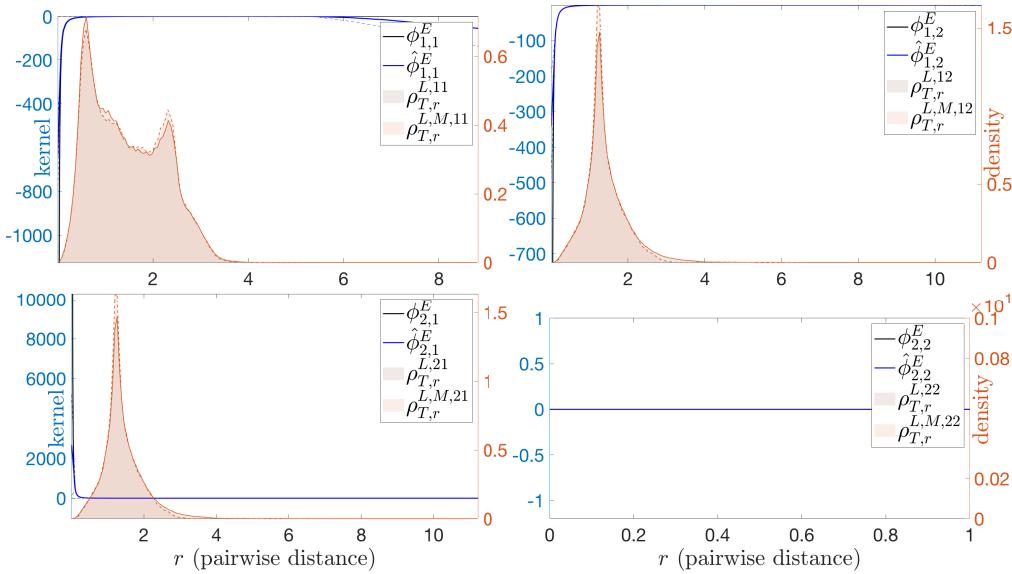


Fig. S11. (PS2<sup>nd</sup>) Interaction kernels learned with  $\text{Unif.}([-\sigma, \sigma])$  multiplicative noise, for  $\sigma = 0.1$  in the observed positions and velocities, with parameters as in Table S12. The estimated kernels are minimally affected, mostly in regions with small  $\rho_T^L$  near 0.

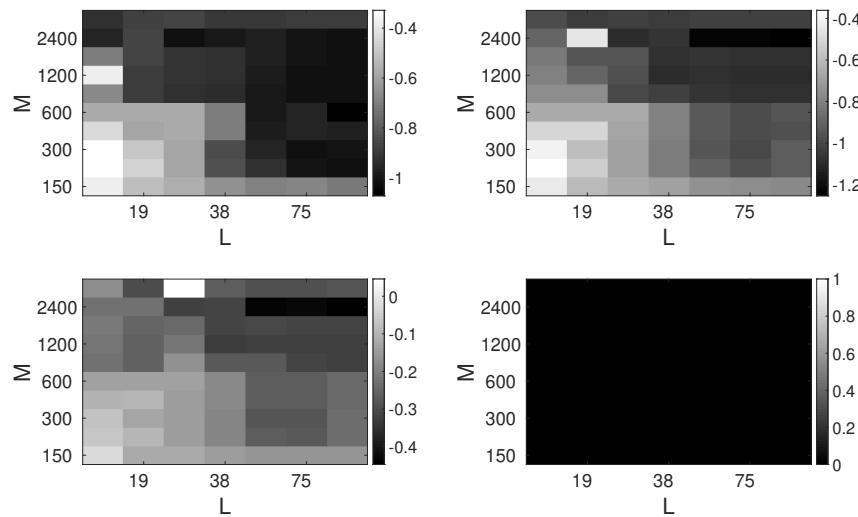


Fig. S12. (PS2) Relative error, in  $\log_{10}$  scale, of  $\phi_{k,k'}^E$  (with  $(k, k')$  increasing lexicographically from top-left to bottom-right) as a function of  $L$  and  $M$ . The error decreases both in  $L$  and  $M$ , in fact roughly in the product  $ML$  (we impute the lack of monotonicity of some of the entries in the plots to the variance in the results). The fourth plot is an identically 0 absolute error, because both  $\phi_{2,2}^E$  and its estimator are identically 0, since there is only one predator. Note  $M \gg 1$  seems to be needed for accurate inference of the interaction kernels, regardless of how large  $L$  is: the trajectories explored for small  $M$  do not explore enough configuration to enable estimation, suggesting that the limit  $M \rightarrow +\infty$  considered in this work is of fundamental importance, at least for non-ergodic systems.

2481 **D. Phototaxis Dynamics.** Second order models have been widely used in describing self-organized human motion (21–23),  
 2482 synthetic agent (robots, drones, etc.) behavior (24–27), and bacteria/cell aggregation and motility (28–31). A step further  
 2483 in accurately model reality is to consider models with responses of agents to their surrounding environment or the spread of  
 2484 emotion among agents within a system. Such phenomena appear in a variety of applications, including modeling of emergency  
 2485 evacuation, crowded pedestrian dynamics, bacteria movement toward certain food sources (28–36). We choose here a system  
 2486 modeling the dynamics of phototactic bacteria towards a fixed light source. This system extends the Cucker-Smale system  
 2487 (9, 37, 38) with an extra auxiliary variable  $\xi_i$  modeling the response (called excitation level) of individual bacteria to the light  
 2488 source. The dynamics is known to lead to flocking (all bacteria moving in the same direction) within a rather short amount  
 2489 time, due to the interaction kernel having a long interaction range and the effect of light entering the dynamics uniformly. This  
 2490 system is within our family of the second order systems, with homogeneous agents and no energy-induced interaction kernel.  
 2491 The alignment-based interaction kernels acting on  $\dot{x}_i$  and  $\xi_i$  are the same:  
 2492  
 2493  
 2494

$$\phi^v(r) = \phi^\xi(r) = (1 + r^2)^{-\frac{1}{4}}.$$

2495  
 2496  
 2497 The non-collective change on  $\dot{x}_i$  is given by  
 2498  
 2499  
 2500

$$F_i^v(\dot{x}_i, \xi_i) = I_0(v_{\text{term}} - \dot{x}_i)(1 - \gamma(\xi_i; \xi_{\text{cr}})),$$

2501 where  $I_0 = 0.1$  is the light intensity,  $v_{\text{term}} = (60, 0)$  is the terminal velocity (light source at infinity),  $\xi_{\text{cr}} = 0.3$  is the critical  
 2502 excitation level (when the light effect activates the bacteria), and  $\gamma(\cdot)$  is the smooth cutoff function  
 2503  
 2504

$$\gamma(\xi; \xi_c) = \begin{cases} 1, & 0 \leq \xi < \xi_c, \\ \frac{1}{2}(\cos(\frac{\pi}{\xi_c}(\xi - \xi_c) + 1), & \xi_c \leq \xi < 2\xi_c, \\ 0, & 2\xi_c \leq \xi. \end{cases}$$

2505 Here  $\xi_c$  is a threshold constant. The non-collective change on  $\xi_i$  is given by  
 2506  
 2507

$$F_i^\xi(\xi_i) = I_0\gamma(\xi_i; \xi_{\text{cp}}),$$

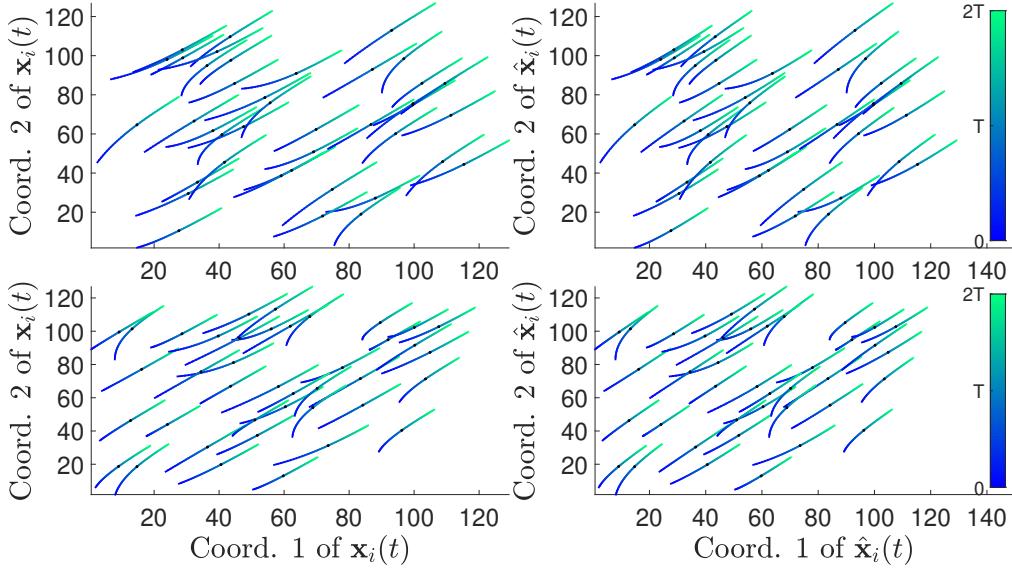
2508 where  $\xi_{\text{cp}} = 0.6$  is the maximum excitation level of light effect on the bacteria. The system parameters are summarized in  
 2509 Table S15.  
 2510  
 2511

**Table S15**

$d$	$M$	$L$	$T$
2	50	200	0.25
$\mu_0^X = \mu_0^{\bar{X}}$	$\mu_0^{\bar{\Xi}}$	$n^v = n^\xi$	$\deg(\psi_{kk'}^A) = \deg(\psi_{kk'}^\xi)$
Unif. on $[0, 100]^2$	Unif. on $[0, 0.001]^2$	400	1

(PT) Parameters for Phototaxis Dynamics

2512  
 2513  
 2514 In the right column of Fig. 5 in the main text, we show the comparison of the learned interaction kernels  $\hat{\phi}^A$  and  $\hat{\phi}^\xi$  versus  
 2515 the true interaction kernels, as well as the comparison of true and learned trajectories over two different sets of initial conditions.  
 2516 We are able to accurately learn the interaction kernels  $\hat{\phi}^A$  and  $\hat{\phi}^\xi$  over the support of  $\rho_T$  when pairwise distance data is  
 2517 abundant. When the pairwise distance data becomes scarce towards the two ends of the interaction interval  $[0, R]$ , we are able  
 2518 to faithfully capture the behavior of  $\phi$  at  $r = 0$ ; the errors are larger near the upper end  $r = R$ , where the data is extremely  
 2519 scarce. Crucially, we recover faithfully the interactions between the agents and their environment. Estimation errors in the  
 2520 appropriate  $L^2(\rho_{T,r,\dot{r}}^L)$ - and  $L^2(\rho_{T,r,\xi}^L)$ -norms are reported in Table S16. A case with noisy observation is also investigated and  
 2521 shown in Fig. S15. Trajectory errors are shown in Table S17. We also compare in Fig. S13 the true and learned trajectories for  
 2522 a corresponding system a dynamics with larger  $N$ .  
 2523  
 2524



**Fig. S13.** (PT) Trajectories  $\mathbf{X}(t)$  and  $\hat{\mathbf{X}}(t)$  obtained with true and learned interaction kernels respectively, for two randomly chosen initial conditions and evolved using the larger number of agents  $N_{\text{new}}$  (governed by the same equations as in the case of  $N$  agents). Trajectory errors are shown in Table S17.

**Table S16**

Rel. Err. for $\hat{\phi}^A$	$9.4 \cdot 10^{-3} \pm 5.2 \cdot 10^{-3}$
Rel. Err. for $\hat{\phi}^\xi$	$8.2 \cdot 10^{-3} \pm 5.0 \cdot 10^{-3}$

(PT) Estimator Errors

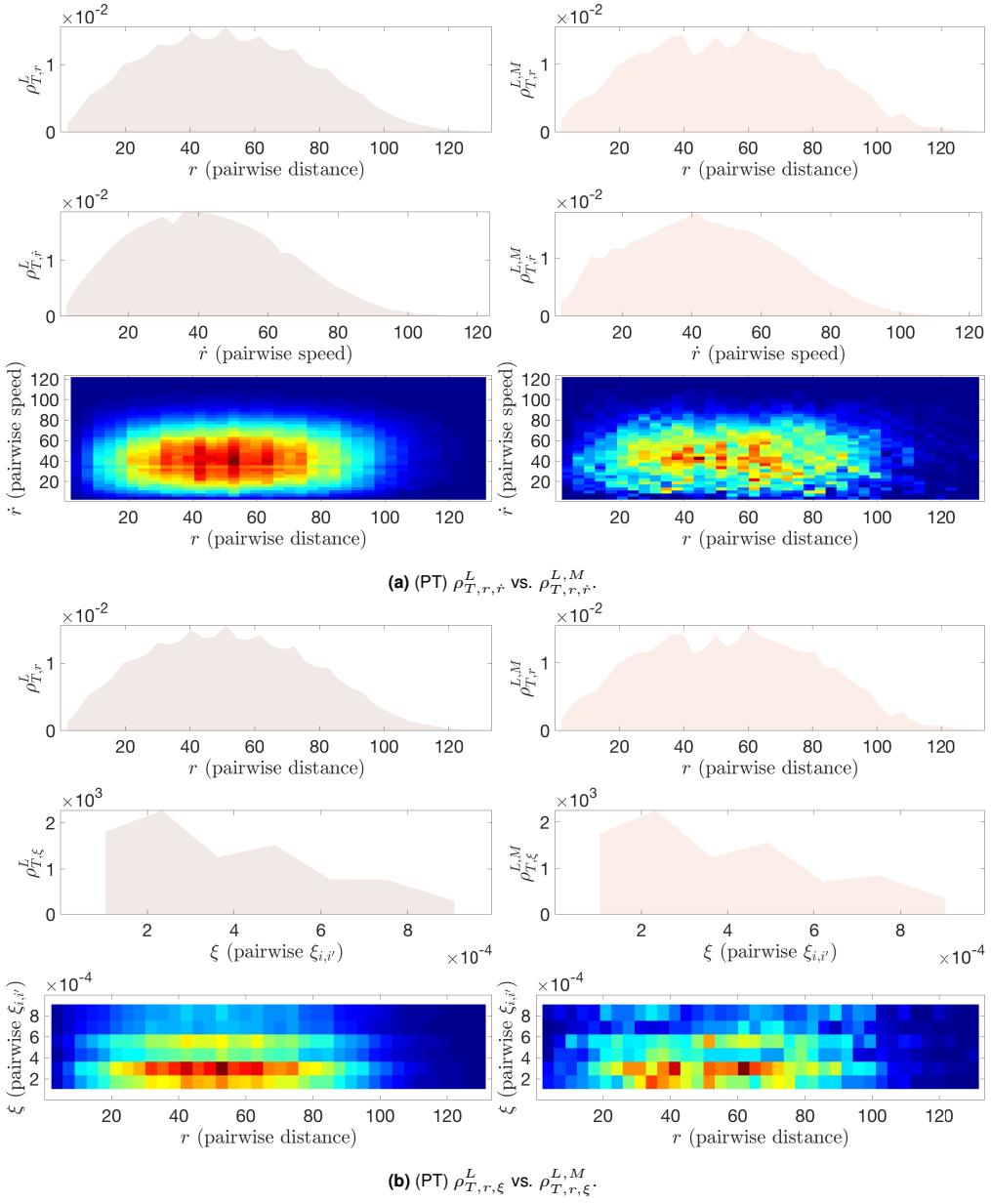
**Table S17**

	[0, T]	[T, T_f]
mean <sub>IC</sub> : Training ICs	$1.6 \cdot 10^{-3} \pm 5.7 \cdot 10^{-5}$	$6.5 \cdot 10^{-3} \pm 9.1 \cdot 10^{-4}$
std <sub>IC</sub> : Training ICs	$3.1 \cdot 10^{-4} \pm 4.8 \cdot 10^{-5}$	$8.1 \cdot 10^{-3} \pm 3.9 \cdot 10^{-3}$
mean <sub>IC</sub> : Random ICs	$1.8 \cdot 10^{-3} \pm 8.0 \cdot 10^{-4}$	$7.3 \cdot 10^{-3} \pm 3.2 \cdot 10^{-3}$
std <sub>IC</sub> : Random ICs	$1.5 \cdot 10^{-3} \pm 3.4 \cdot 10^{-3}$	$1.1 \cdot 10^{-2} \pm 1.2 \cdot 10^{-2}$
mean <sub>IC</sub> : Larger N	$4.2 \cdot 10^{-3} \pm 1.6 \cdot 10^{-3}$	$8.4 \cdot 10^{-3} \pm 3.8 \cdot 10^{-3}$
std <sub>IC</sub> : Larger N	$2.9 \cdot 10^{-3} \pm 3.0 \cdot 10^{-3}$	$7.9 \cdot 10^{-3} \pm 7.0 \cdot 10^{-3}$

(PT) Trajectory Errors

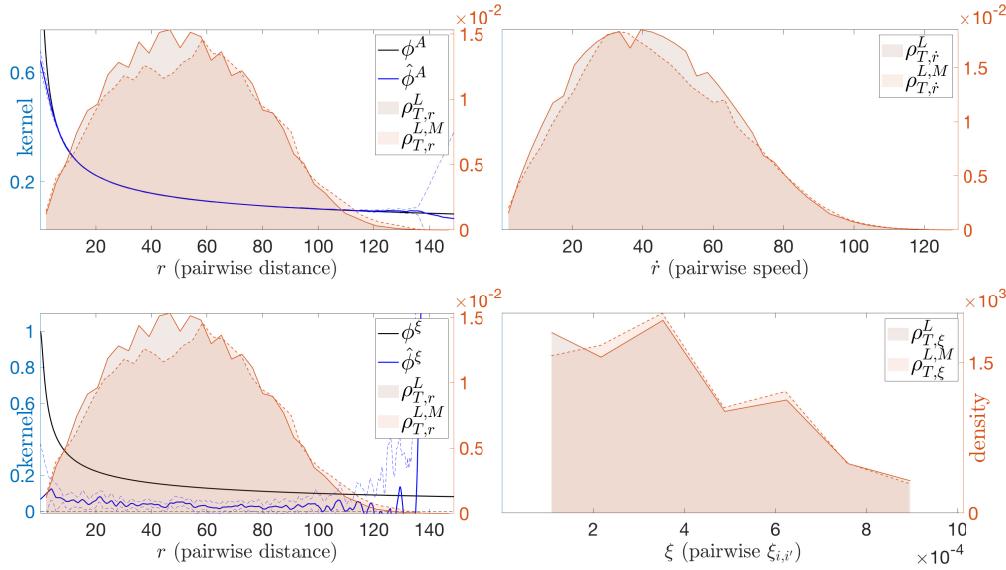
x

Finally we display, in Fig. S14a and S14b, the two joint distributions  $\rho_{T,r,\hat{r}}^L$  and  $\rho_{T,r,\xi}^L$ , used to define the appropriate  $L^2$ -norms for measuring the performance of  $\hat{\phi}^A$  and  $\hat{\phi}^\xi$ . We also calculated the  $\ell^1$  distance between the joint distribution  $\rho_{T,r,\hat{r}}^L$  and the product of its marginals, and it is  $1 \cdot 10^{-1}$ . For the  $\ell^1$  distance between  $\rho_{T,r,\xi}^L$  and the product of its marginals, it is  $7 \cdot 10^{-2}$ . For the empirical distributions (over 10 learning trials), the  $\ell^1$  distance for  $\rho_{T,r,\hat{r}}^{L,M}$  and the product of its marginal is  $7 \cdot 10^{-1} \pm 1 \cdot 10^{-2}$ ; whereas the  $\ell^1$  distance of  $\rho_{T,r,\xi}^{L,M}$  to the product of its marginals is  $3.7 \cdot 10^{-1} \pm 7 \cdot 10^{-3}$ .



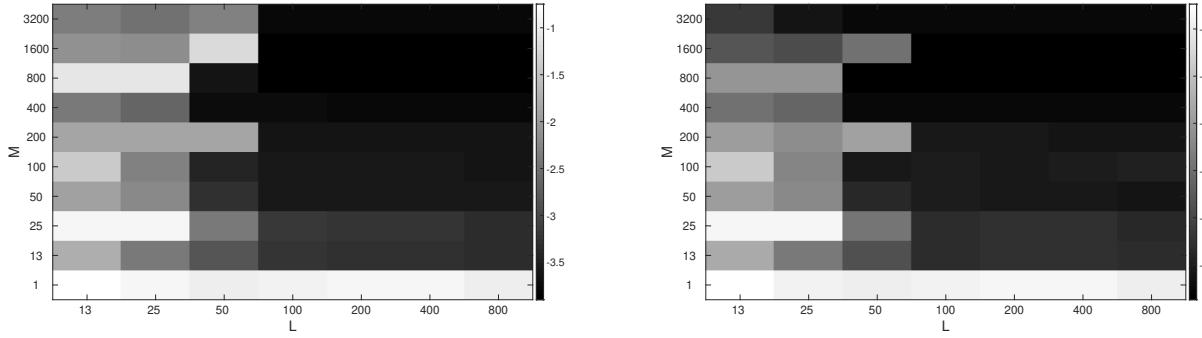
**Fig. S14.** (PT) Density plots for the various  $\rho_T^L$  measures.

2729  
2730  
2731  
2732  
2733  
2734  
2735  
2736  
2737  
2738  
2739  
2740  
2741  
2742  
2743  
2744  
2745  
2746  
2747  
2748  
2749  
2750  
2751  
2752  
2753  
2754  
2755  
2756  
2757  
2758  
2759  
2760  
2761  
2762  
2763  
2764  
2765  
2766  
2767  
2768  
2769  
2770  
2771  
2772  
2773  
2774  
2775  
2776  
2777  
2778  
2779  
2780  
2781  
2782  
2783  
2784  
2785  
2786  
2787  
2788  
2789  
2790  
2791  
2792  
2793  
2794  
2795  
2796  
2797  
2798  
2799  
2800  
2801  
2802  
2803  
2804  
2805  
2806  
2807  
2808  
2809  
2810  
2811  
2812  
2813  
2814  
2815  
2816  
2817  
2818  
2819  
2820  
2821  
2822  
2823  
2824  
2825  
2826  
2827  
2828  
2829  
2830  
2831  
2832  
2833  
2834  
2835  
2836  
2837  
2838  
2839  
2840  
2841  
2842  
2843  
2844  
2845  
2846  
2847  
2848  
2849  
2850  
2851  
2852



**Fig. S15.** (PT) Interaction kernels learned from noisy observations of positions and velocities. The noises are multiplicative,  $\text{Unif}([-\sigma, \sigma])$  with  $\sigma = 0.1$  and with other parameters as in Table S15. The estimated kernel for associated with  $\dot{x}_i$  is minimally affected, mostly in regions with small  $\rho_T^L$ ; the additive noise is on a scale far greater than that on  $\xi_i$  hence severely affects the learning result on the interaction kernel on  $\xi_i$ .

Figure S16 shows the behavior of the error of the estimators as both  $L$  and  $M$  are increased.



**Fig. S16.** (PT) Relative error, in  $\log_{10}$  scale, of  $\hat{\phi}^A$  (left) and  $\hat{\phi}^\xi$  (right) as a function of  $L$  and  $M$ . The error decreases both in  $L$  and  $M$ , in fact roughly in the product  $ML$ . The fourth plot is an identically 0 absolute error, because both  $\phi_{2,2}^E$  and its estimator are identically 0, since there is only one predator. Note  $M \gg 1$  seems to be needed for accurate inference of the interaction kernels, regardless of how large  $L$  is: the trajectories explored for small  $M$  do not explore enough configuration to enable estimation, suggesting that the limit  $M \rightarrow +\infty$  considered in this work is of fundamental importance, at least for non-ergodic systems.

**E. Model Selection.** Our learning approach can be used to identify the model of the system from the observation data. We consider here two different scenarios of model selection: one is identifying the type – energy-based vs. alignment-based – of interaction kernels from a second order system driven by only one type of interaction kernel; the other is to identify the order of the system from a heterogeneous dynamics.

*Model Selection: energy-based vs. alignment-based interactions.* We consider a special case of the second order homogeneous agent dynamics, given as either

$$\ddot{x}_i = \sum_{i'=1}^N \frac{1}{N} \phi^E(r_{ii'}) \mathbf{r}_{ii'} \quad \text{or} \quad \ddot{x}_i = \sum_{i'=1}^N \frac{1}{N} \phi^A(r_{ii'}) \dot{\mathbf{r}}_{ii'},$$

with the (unknown) interaction kernels defined as

$$\phi^E(r) = 2 - \frac{1}{r^2} \quad \text{and} \quad \phi^A(r) = \frac{1}{(1+r^2)^{0.25}}.$$

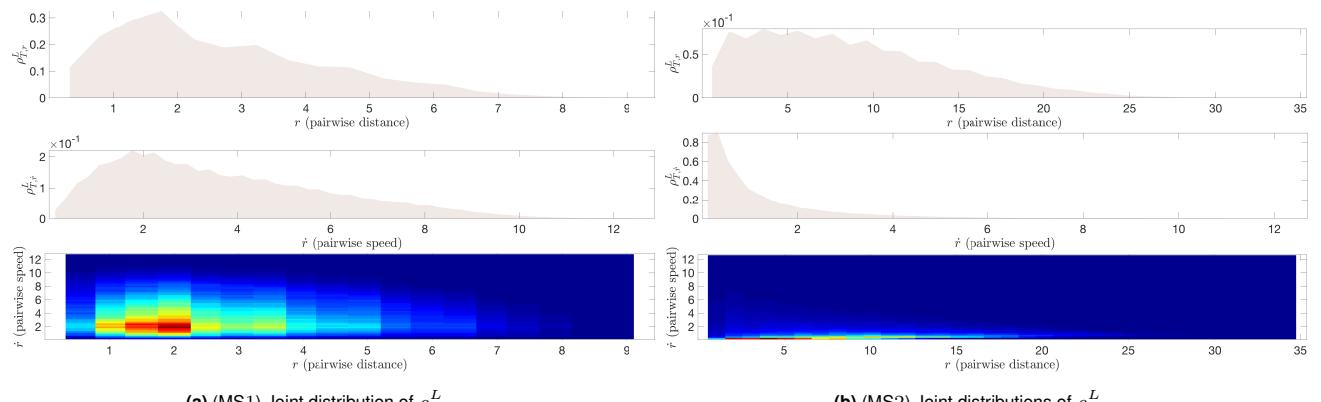
The system parameters are given in Table S18.

**Table S18**

2977	$d$	$M$	$L$	$T$	$\mu_0^X$	$\mu_0^X$	$n^E = n^A$	$\deg(\psi^A) = \deg(\psi^\xi)$	3039
2978	2	200	200	10	Unif. on ring [0.5, 1]	$\mathcal{U}([0, 10]^2)$	800	1	3040
(MS1 and 2) Test Parameters									3041

2981  
2982  
2983

2984 Given the observation data from either system ( $\phi^E$ - or  $\phi^A$ -driven), we proceed to learn the interaction kernels as usual,  
2985 i.e. as if the dynamics were generated with both energy-based and alignment-based interaction kernels present. Results are  
2986 shown in Fig. 7 in the main text. The two sub-figures on the left show the learned interaction kernels  $\hat{\phi}^E$  and  $\hat{\phi}^A$  from a purely  
2987 energy-based system:  $\hat{\phi}^A$  is small in the appropriate norm, while  $\hat{\phi}^E$  is large (and a good approximation to  $\phi^E$ ): the estimators  
2988 can therefore detect this is an energy-driven system. In the two sub-figures on the right, we display the analogous results  
2989 corresponding to learning the interaction kernels for an alignment-based system. We obtain (almost) 0 for the norm of  $\hat{\phi}^E$ . The  
2990 reason why the  $L^2(\rho_{T,r,\dot{r}}^L)$  norm of  $\hat{\phi}^A$  (from the first case) is not as close to 0 as the  $L^2(\rho_{T,r,\dot{r}}^L)$  norm of the  $\hat{\phi}^E$  (from the second  
2991 case) lies in the difference in the joint distribution of the two cases, see Figures S17a and S17b. To further investigate the  
2992 properties of the joint distributions (and also to differentiate the two dynamics), we calculated the  $\ell^1$  distance of the respective  
2993 joint distributions to the product and their marginals. For MS1, the  $\ell^1$  distance (over 10 learning trials) between the joint  
2994 distribution  $\rho_{T,r,\dot{r}}^{L,M}$  and the product of its marginals is  $1.3 \cdot 10^{-1} \pm 3.8 \cdot 10^{-3}$ . For MS2, the  $\ell^1$  distance (over 10 learning trials)  
2995 between the joint distribution  $\rho_{T,r,\dot{r}}^{L,M}$  and the product of its marginals is  $4.6 \cdot 10^{-1} \pm 3.4 \cdot 10^{-3}$ .



**Fig. S17.** (MS1 and 2) Density plots for the various  $\rho_T^L$  measures.

3011 *Model Selection: first order vs. second order.* We consider two different heterogeneous agent systems, one first order and one  
3012 second order, with the order of the system unknown to the estimator. The observations are in the time interval  $[0, T]$ , and in  
3013 this case  $T_f = T$ . We first consider the first order heterogeneous agent system

$$\dot{\mathbf{x}}_i = \sum_{i'=1}^N \frac{1}{N_{k_i k_{i'}}} \phi_{k_i k_{i'}}(r_{ii'}) \mathbf{r}_{ii'},$$

3014 with

$$\phi_{1,1}(r) = 1 - r^{-2}, \quad \phi_{1,2}(r) = -2r^{-2}, \quad \phi_{2,1}(r) = 3.5r^{-3}, \quad \phi_{2,2}(r) \equiv 0,$$

3015 and the type information setup similar to that of the Predator-Swarm first order system (detailed in Sec.3C). For the second  
3016 scenario, we consider the data generated by the following second order heterogeneous agent dynamics,

$$\ddot{\mathbf{x}}_i = -\dot{\mathbf{x}}_i + \sum_{i'=1}^N \frac{1}{N_{k_i k_{i'}}} \phi_{k_i k_{i'}}^E(r_{ii'}) \mathbf{r}_{ii'},$$

3017 with

$$\phi_{1,1}(r) = 1 - r^{-2}, \quad \phi_{1,2}(r) = -r^{-2}, \quad \phi_{2,1}(r) = 1.5r^{-2.5}, \quad \phi_{2,2}(r) \equiv 0,$$

3018 and the type information setup similar to that of the Predator-Swarm second order system (details shown in Sec.3C). The  
3019 parameters for both systems are given in Tables S19 and S20.

3101  
3102  
3103  
3104  
3105  
3106  
3107  
3108  
3109  
3110  
3111  
3112  
3113  
3114  
3115  
3116  
3117  
3118  
3119  
3120  
3121  
3122  
3123  
3124  
3125  
3126  
3127  
3128  
3129  
3130  
3131  
3132  
3133  
3134  
3135  
3136  
3137  
3138  
3139  
3140  
3141  
3142  
3143  
3144  
3145  
3146  
3147  
3148  
3149  
3150  
3151  
3152  
3153  
3154  
3155  
3156  
3157  
3158  
3159  
3160  
3161  
3162  
**Table S19**  

$d$	$M$	$L$	$T$
2	250	250	1
$n$	Deg( $\psi_{kk'}$ )	Prey $\mu_0^X$	Pred. $\mu_0^X$
[298, 150; 150, 2]	[1, 1; 1, 0]	Unif. on ring [0.5, 1.5]	Unif. on disk at 0.1

(MS3) Test Parameters

3163  
3164  
3165  
3166  
3167  
3168  
3169  
3170  
3171  
3172  
3173  
3174  
3175  
3176  
3177  
3178  
3179  
3180  
3181  
3182  
3183  
3184  
3185  
3186  
3187  
3188  
3189  
3190  
3191  
3192  
3193  
3194  
3195  
3196  
3197  
3198  
3199  
3200  
3201  
3202  
3203  
3204  
3205  
3206  
3207  
3208  
3209  
3210  
3211  
3212  
3213  
3214  
3215  
3216  
3217  
3218  
3219  
3220  
3221  
3222  
3223  
3224

3111  
3112  
3113  
3114  
3115  
3116  
3117  
3118  
3119  
3120  
3121  
3122  
3123  
3124  
3125  
3126  
3127  
3128  
3129  
3130  
3131  
3132  
3133  
3134  
3135  
3136  
3137  
3138  
3139  
3140  
3141  
3142  
3143  
3144  
3145  
3146  
3147  
**Table S20**  

$d$	$M$	$L$	$T$
2	250	250	1
$n$	deg( $\psi_{kk'}^E$ )	Prey $\mu_0^X$	Pred. $\mu_0^X$
[298, 150; 150, 2]	[1, 1; 1, 0]	$\mathcal{U}([0.1, 1]^2)$	$\mathcal{U}([0, 0.07]^2)$

(MS4) Test Parameters

With the order of the ODE system and the interaction kernels being the missing information, we construct estimators for the interaction kernels in two ways: first assuming a first order system, then assuming a second order system (without non-collective forcing). We then generate predicted trajectories using the learned interaction kernels, and the same initial conditions as in the training data. Next, we calculate the trajectory max-in-time error, obtaining the results in Table 1 of the main text (shown as the mean of the trajectory error plus or minus standard deviation of the error over 10 runs). As indicated by the trajectory error statistics, the predicted trajectories with smaller error indicate the correct order of the true underlying system in both cases. Details on the statistics of the trajectory errors are reported in Tables S21 and S22. In each, the column with smaller values (within both mean and standard deviation of the trajectory errors) corresponds the correct order of the system.

3148  
3149  
3150  
3151  
3152  
3153  
3154  
3155  
3156  
3157  
3158  
3159  
3160  
3161  
3162  
**Table S21**  

	Learned as 1 <sup>st</sup> order	Learned as 2 <sup>nd</sup> order
mean <sub>IC</sub>	$9.5 \cdot 10^{-3} \pm 2 \cdot 10^{-3}$	$3.9 \pm 8$
std <sub>IC</sub>	$1.8 \cdot 10^{-2} \pm 1.1 \cdot 10^{-2}$	$48 \pm 1 \cdot 10^2$

(MS3) Trajectory Errors

3143  
3144  
3145  
3146  
3147  
3148  
3149  
3150  
3151  
3152  
3153  
3154  
3155  
3156  
3157  
3158  
3159  
3160  
3161  
3162  
**Table S22**  

	Learned as 1 <sup>st</sup> order	Learned as 2 <sup>nd</sup> order
mean <sub>IC</sub>	$1.6 \pm 1 \cdot 10^{-1}$	$1.3 \cdot 10^{-1} \pm 3 \cdot 10^{-2}$
std <sub>IC</sub>	$9.4 \cdot 10^{-1} \pm 2 \cdot 10^{-1}$	$2.0 \cdot 10^{-1} \pm 5 \cdot 10^{-2}$

(MS4) Trajectory Errors

## References

1. Cucker F, Smale S (2002) On the mathematical foundations of learning. *Bull Amer Math Soc* 39(1):1–49.
2. Binev P, Cohen A, Dahmen W, DeVore R, Temlyakov V (2005) Universal algorithms for learning theory part i: piecewise constant functions. *J Mach Learn Res* 6(Sep):1297–1321.
3. DeVore R, Kerkyacharian G, Picard D, Temlyakov V (2006) Approximation methods for supervised learning. *Found Comput Math* 6(1):3–58.
4. Bongini M, Fornasier M, Hansen M, Maggioni M (2017) Inferring interaction rules from observations of evolutive systems I: The variational approach. *Math Mod Methods Appl Sci* 27(05):909–951.
5. Carrillo JA, D’Orsogna MR, Panferov V (2009) Double Milling in self-propelled swarms from kinetic theory. *Kinet Relat Mod* 2(2):363 – 378.
6. Chuang Y, D’Orsogna M, Marthaler D, Bertozzi A, Chayes L (2007) State transition and the continuum limit for the 2D interacting, self-propelled particle system. *Physica D* 232:33 – 47.
7. Cristiani E, Piccoli B, Tosin A (2010) Modeling self-organization in pedestrians and animal groups from macroscopic and microscopic viewpoints in *Mathematical Modeling of Collective Behavior in Socio-Economic and Life Sciences, Modeling and Simulation in Science, Engineering and Technology*, eds. Naldi G, Pareschi L, Toscani G, Bellomo N. (Springer, Birkhäuser Boston), pp. 337 – 364.

- 3225 8. Couzin I, Franks N (2002) Self-organized lane formation and optimized traffic flow in army ants. *Proc R Soc Lond B* 3287  
 3226 270:139 – 146. 3288
- 3227 9. Cucker F, Smale S (2007) Emergent behavior in flocks. *IEEE Trans Automat Contr* 52(5):852. 3289
- 3228 10. Niwa H (1994) Self-organizing dynamic model of fish schooling. *J Theor Biol* 171:123 – 136. 3290
- 3229 11. Parrish JK, Edelstein-Keshet L (1999) Complexity, pattern, and evolutionary trade-offs in animal aggregation. *Science* 3291  
 3230 284:99 – 101. 3292
- 3231 12. Parrish J, Viscido S, Grunbaum D (2002) Self-organized fish schools: An examination of emergent properties. *Biol Bull* 3293  
 3232 202:296 – 305. 3294
- 3233 13. Romey W (1996) Individual differences make a difference in the trajectories of simulated schools of fish. *Ecol Model* 92:65 – 3295  
 3234 77. 3296
- 3235 14. Toner J, Tu Y (1995) Long-range order in a two-dimensional dynamical xy model: How birds fly together. *Phys Rev Lett* 3297  
 3236 75:4326 – 4329. 3298
- 3237 15. Yates C, et al. (2009) Inherent noise can facilitate coherence in collective swarm motion. *Proc Natl Acad Sci USA* 106:5464 – 3299  
 3238 5469. 3300
- 3239 16. Escobedo R, Muro C, Spector L, Coppinger RP (2014) Group size, individual role differentiation and effectiveness of 3301  
 3240 cooperation in a homogeneous group of hunters. *J R Soc Interface* 11:20140204. 3302
- 3241 17. Cohn H, Kumar A (2009) Algorithmic design of self-assembling structures. *Proc Natl Acad Sci USA* 106:9570 – 9575. 3303
- 3242 18. Nowak MA (2006) Five rules for the evolution of cooperation. *Science* 314:1560 – 1563. 3304
- 3243 19. Fryxell JM, Mosser A, Sinclair ARE, Packer C (2007) Group formation stabilizes predator-prey dynamics. *Nature* 449:1041 3305  
 3244 – 1043. 3306
- 3245 20. Cucker F, Dong JG (2014) A conditional, collision-avoiding, model for swarming. *Discrete Continuous Dyn Syst* 43(3):1009 3307  
 3246 – 1020. 3308
- 3247 21. Cristiani E, Piccoli B, Tosin A (2011) Multiscale modeling of granular flows with application to crowd dynamics. *Multi 3309  
 3248 Model Simul* 9(1):155 – 182. 3310
- 3249 22. Cucker F, Smale S, Zhou D (2004) Modeling language evolution. *Found Comput Math* 4(5):315 – 343. 3311
- 3250 23. Short MB, et al. (2008) A statistical model of criminal behavior. *Math Models Methods Appl Sci* 18(suppl.):1249 – 1267. 3312
- 3251 24. Chuang Y, Huang Y, D'Orsogna M, Bertozzi A (2007) Multi-vehicle flocking: scalability of cooperative control algorithms 3313  
 3252 using pairwise potentials. *IEEE Intern Conf Robotics and Automation* pp. 2292 – 2299. 3314
- 3253 25. Leonard N, Fiorelli E (2001) Virtual leaders, artificial potentials and coordinated control of groups. *Proc 40<sup>th</sup> IEEE Conf 3315  
 3254 Decision Contr* pp. 2968 – 2973. 3316
- 3255 26. Pera L, Gómez G, Elosegui P (2009) Extension of the Cucker-Smale control law to space flight formations. *J Guid Control 3317  
 3256 Dyn* 32:527 – 537. 3318
- 3257 27. Sugawara K, Sano. M (1997) Cooperative acceleration of task performance: Foraging behavior of interacting multi-robots 3319  
 3258 system. *Physica D* 100:343 – 354. 3320
- 3259 28. Camazine S, et al. (2001) *Self-organization in Biological Systems*, Princeton studies in complexity. (Princeton University 3321  
 3260 Press, Princeton). 3322
- 3261 29. Evelyn FK, Lee AS (1970) Initiation of slime mold aggregation viewed as an instability. *J Theor Biol* 26 3:399–415. 3323
- 3262 30. Koch AL, White D (1998) The social lifestyle of myxobacteria. *BioEssays* 20(12):1030–1038. 3324
- 3263 31. Perthame B (2007) *Transport Equations in Biology*, Frontiers in Mathematics. (Birkhäuser Basel). 3325
- 3264 32. Moussaid M, Helbing D, Theraulaz G (2011) How simple rules determine pedestrian behavior and crowd disasters. *Proc 3326  
 3265 Natl Acad Sci USA* 108(17):6884 – 6888. 3327
- 3266 33. Durupinar F, Gudukbar U, Aman A, Badler NI (2015) Psychological Parameters for Crowd Simulation: From Audiences 3328  
 3267 to Mobs. *IEEE Trans Vis Comput Graph* 21:1 – 15. 3329
- 3268 34. Bosse T, Duell R, Memon ZA, Treur J, van der Wal CN (2009) A multi-agent model for mutual absorption of emotions in 3330  
 3269 *European council on modeling and simulation*, ECMS 2009. 3331
- 3270 35. Bosse T, Hoogendoorn M, Klein MCA, Treur J, van der Wal CN (2011) Agent-based analysis of patterns in crowd behaviour 3332  
 3271 involving contagion of mental states in *Modern Approaches in Applied Intelligence: 24th International Conference on 3333  
 3272 Industrial Engineering and Other Applications of Applied Intelligent Systems, IEA/AIE 2011, Syracuse, NY, USA, June 3334  
 3273 28 – July 1, 2011, Proceedings, Part II*, eds. Mehrotra KG, Mohan CK, Oh JC, Varshney PK, Ali M. (Springer Berlin 3335  
 3274 Heidelberg, Berlin, Heidelberg), pp. 566–577. 3336
- 3275 36. Lin J, Luckas TA (2015) A particle swarm optimization model of emergency airplane evacuation with emotion. *Net Het 3337  
 3276 Media* 10:631 – 646. 3338
- 3277 37. Cucker F, Smale S (2007) On the mathematics of emergence. *Jpn J Math* 2(1):197 – 227. 3339
- 3278 38. Ha SY, Ha T, Kim JH (2010) Emergent behavior of a Cucker-Smale type particle model with nonlinear velocity couplings. 3340  
 3279 *IEEE Trans Automat Contr* 55(7):1679 – 1683. 3341
- 3280 3342
- 3281 3343
- 3282 3344
- 3283 3345
- 3284 3346
- 3285 3347
- 3286 3348