

Supporting Information: Time Series Modeling of Varying
Complexity Captures Subdiffusive Solute Dynamics and Predicts
Long Timescale Behavior in Nanoscale Pores.

Benjamin J. Coscia Michael R. Shirts

January 8, 2020

Contents

S1	Setup and analysis scripts.....	S2
S2	Solute Equilibration	S3
S3	Estimating the Hurst Parameter.....	S4
S4	Simulating Fractional Lévy Motion	S5
	S4.1 Truncated Lévy stable hop distributions	S5
	S4.2 Achieving the right correlation structure.....	S6
S5	Verifying Markovianity	S7
S6	Derivation of Passage Time Distributions.....	S10
S7	Solute hopping and trapping behavior.....	S11
S8	AD model MSD Predictions with Pure Power Law Dwell Times	S12
S9	Stationarity of Solute Trajectories.....	S13
S10	Tabular Anomalous Diffusion Parameters	S16
S11	Tabular MSDDM parameters	S17
S12	Explanation of Under-Estimated MSDDM β Values	S21

S1 Setup and analysis scripts

All python and bash scripts used to set up systems and conduct post-simulation trajectory analysis are available online at https://github.com/shirtsgroup/LLC_Membranes. Documentation for the `LLC_Membranes` repository is available at <https://llc-membranes.readthedocs.io/en/latest/>. Table S1 provides more detail about specific scripts used for each type of analysis performed in the main text.

Script Name	Section	Description
<code>/setup/parameterize.py</code>	2.1	Parameterize liquid crystal monomers and solutes with GAFF
<code>/setup/build.py</code>	2.1	Build simulation unit cell
<code>/setup/place_solutes_pores.py</code>	2.1	Place equispaced solutes in the pore centers of a unit cell
<code>/setup/equil.py</code>	2.1	Equilibrate unit cell and run production simulation
<code>/analysis/solute_partitioning.py</code>	2.1	Determine time evolution of partition of solutes between pores and tails
<code>/timeseries/msd.py</code>	2.2	Calculate the mean squared displacement of solutes
<code>/analysis/sfbm_parameters.py</code>	2.2	Get subordinated fractional Brownian motion parameters by fitting to a solute's dwell and hop length distributions and positional autocorrelation function.
<code>/timeseries/ctrwsim.py</code>	2.2	Generate realizations of a continuous time random walk with the user's choice of dwell and hop distributions
<code>/timeseries/forecast_ctrw.py</code>	2.2	Combines classes from <code>sfbm_parameters.py</code> and <code>ctrwsim.py</code> to parameterize and predict MSD in one shot.
<code>/analysis/markov_state_dependent_dynamics.py</code>	2.3	Identify frame-by-frame state of each solute, construct a transition matrix and simulate realizations of the MSDDM model.
<code>/timeseries/mfpt_pore.py</code>	2.4	Simulate mean first passage time distributions using the AD or MSDDM model.

Table S1: The first column provides the names of the python scripts available in the `LLC_Membranes` GitHub repository that were used for system setup and post-simulation trajectory analysis. Paths preceding script names are relative to the `LLC_Membranes/LLC_Membranes` directory. The second column lists the section in the main text where the output or usage of the script is first described. The third column gives a brief description of the purpose of each script.

S2 Solute Equilibration

We collected all data used for model generation after the solutes were equilibrated. We assumed a solute to be equilibrated when the partition of solutes in and out of the pore region stopped changing. The pore region is defined as within 0.75 nm of the pore center. We've plotted the partition versus time in Figure S1 and indicated the chosen equilibration time points.

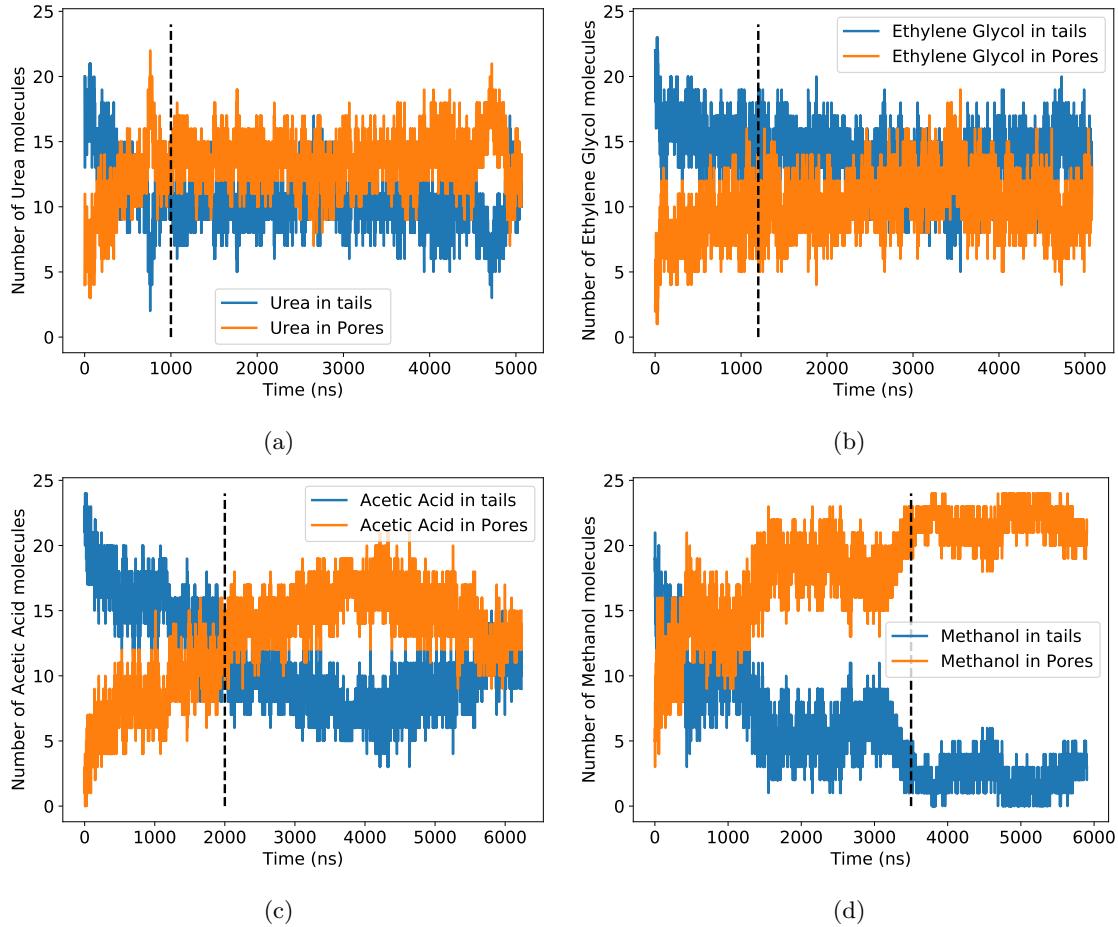


Figure S1: We considered a system to be equilibrated when the partition of solutes between the tails and pore plateaued. Our chosen equilibration point for each solute is indicated by the vertical black dashed line. (a) Urea equilibrates the fastest, after 1000 ns. (b) Ethylene glycol equilibrates after 1200 ns (c) The partition of acetic acid appears oscillate slowly. We considered it to be equilibrated after 2000 ns. (d) We considered methanol to be equilibrated after 3500 ns. Methanol nearly completely partitions into the tails.

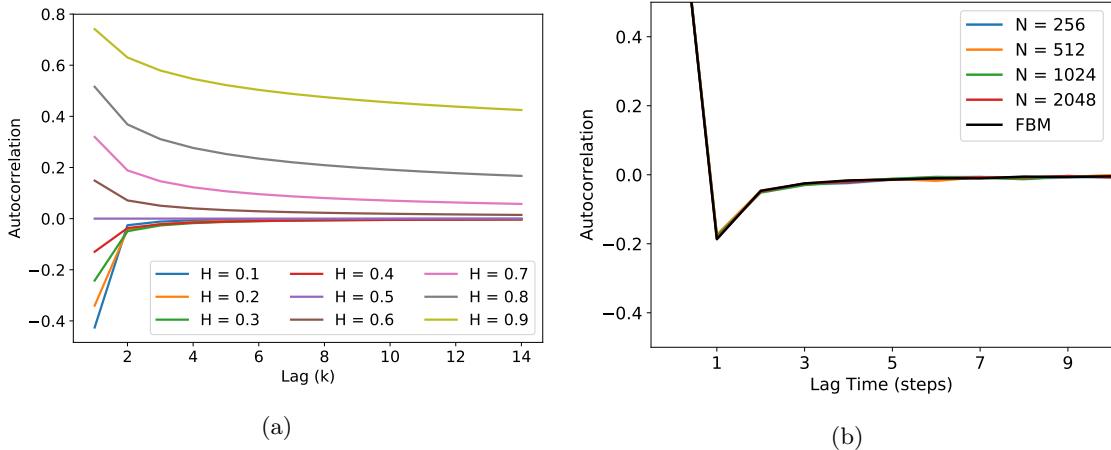


Figure S2: (a) The analytical autocorrelation function of FBM decays to zero faster when $H < 0.5$ compared to when $H > 0.5$. (b) The autocorrelation function of an FLM process does not change with increasing sequence length (N). It shares the same autocorrelation function as fractional Brownian motion (FBM). All sequences used to make this plot were generated using $H=0.35$ and, for FLM, $\alpha=1.4$.

S3 Estimating the Hurst Parameter

We chose to estimate the Hurst parameter, H by a least squares fit to the analytical autocorrelation function for fractional Brownian motion (the variance-normalized version of Equation 6 in the main text):

$$\gamma(k) = \frac{1}{2} \left[|k-1|^{2H} - 2|k|^{2H} + |k+1|^{2H} \right] \quad (1)$$

In Figure S2a, we plotted Equation 1 for different values of H . When $H > 0.5$, Equation 1 decays slowly to zero meaning one needs to study large time lags with high frequency in order to accurately estimate H from the data. Fortunately, all of our solutes show anti-correlated motion, so most of the information in Equation 1 is contained within the first few lags.

The autocovariance function of fractional Lévy motion is different from fractional Brownian motion (see Equations 6 and 8 of the main text), but their autocorrelation structures are the same. The autocovariance function of FLM is dependent on the expected value of squared draws from the underlying Lévy distribution, $E[L(1)^2]$. This is effectively the distribution's variance, which is undefined for most Lévy stable distributions due to their heavy tails. As a consequence, one should expect $E[L(1)^2]$ to grow as more samples are drawn from the distribution with the autocovariance function responding accordingly. However, we are only interested in the autocorrelation function. In order to predict the Hurst parameter from the autocorrelation function, we must show that it has a well-defined structure and is independent of the coefficient in Equation 8 of the main text. In Figure S2b, we plot the average autocorrelation function from an FLM process with an increasing number of observations per generated sequence. For all simulations we set $H=0.35$ and $\alpha=1.4$. The variance-normalized autocovariance function, i.e. the autocorrelation function, does not change with increasing sequence length. Additionally, the autocorrelation function of FBM, with the same H , is the same. Therefore we are confident that we can use the same Hurst parameter as an input to both FBM and FLM simulations.

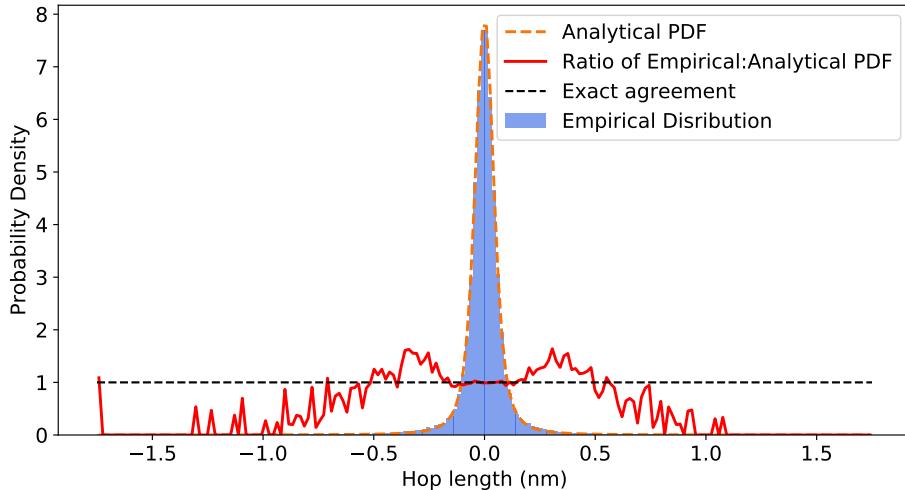


Figure S3: The ratio between the empirical and maximum likelihood theoretical distribution quantifies the quality of fit as function of hop length. The fit is near-perfect close to the mean. Intermediate hop lengths are over sampled, and the tails are under sampled. We used this type of plot to determine the appropriate place to truncate the Lévy stable distributions.

S4 Simulating Fractional Lévy Motion

S4.1 Truncated Lévy stable hop distributions

Determining where to truncate the hop distribution: A pure Lévy stable distribution has heavy tails which can lead to arbitrarily long hop lengths. Our distribution of hop lengths fits well to a Lévy distribution near the mean, but under samples the tails. In Figure S4 we compare the empirically measured transition emission distribution of the MSDDM for urea to its maximum likelihood fit to a Lévy stable distribution. The ratio between the two distributions at each bin is nearly 1 close to the center, indicating a near-perfect fit, larger than 1 slightly further from the center, suggesting that we slightly over sample intermediate hop lengths, and below 1 far from the center, indicating under sampling of extremely long hop lengths. Based on the plot, we chose a cut-off of 1 nm in order to compensate for over sampled intermediate hop lengths. We chose the same cut-off for all solutes.

Generating FLM realizations from a truncated Lévy distribution: To generate realizations from an uncorrelated truncated Lévy process, one would randomly sample from the base distribution and replace values that are too large with new random samples from the base distribution, repeating the process until all samples are under the desired cut-off.

This procedure is complicated by the correlation structure of FLM. At a high level, Stoev and Taqqu use Riemann-sum approximations of the stochastic integrals defining FLM in order to generate realizations. [1] They do this efficiently with the help of Fast Fourier Transforms. In practice, this requires one to Fourier transform a zero-padded vector of random samples drawn from the appropriate Lévy stable distribution, multiply the vector in Fourier space by a kernel function and invert back to real space. The end result is a correlated vector of fractional Lévy noise.

We are unaware of a technique for simulating truncated FLM, therefore we devised our own based on the above discussion. If one is to truncate an FLM process, one can apply the simple procedure above for drawing uncorrelated values from the marginal Lévy stable distribution, *but*, after adding correlation, the maximum drawn value is typically lower than the limit set by the user. Additionally, the shape of the distribution itself changes. Therefore, we created a database meant to correct the input truncation parameter (the maximum desired draw). The database returns the value of the truncation parameter that will properly truncate the output marginal distribution based on H , α and σ (the width parameter). Figure S4 shows the result of applying our correction. Note that generating this database requires a significant amount of simulation and still likely doesn't perfectly correct the truncation parameter. The output leads to a somewhat fuzzy, rather

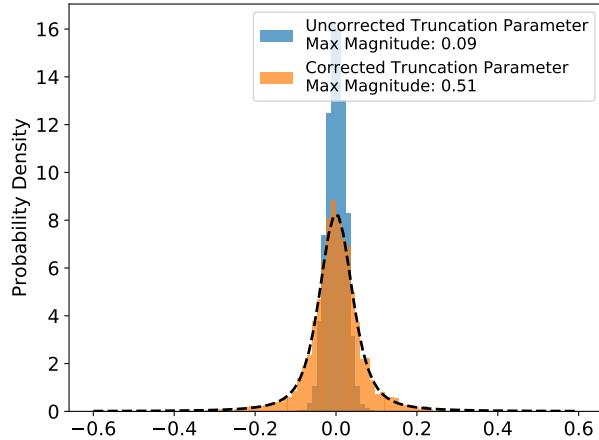


Figure S4: We can accurately truncate the marginal distribution of FLM innovations by applying a correction to the input truncation parameter. We generated FLM sequences and truncated the initial Lévy stable distribution (before Fourier transforming) at a value of 0.5. After correlation structure is added, the width of the distribution of fractional Lévy noise decreases significantly. We corrected the input truncation parameter with our database resulting in a distribution close to the theoretical distribution (black dashed line) with a maximum value close to 0.5.

than abrupt, cut-off of the output distribution. This is likely beneficial since we observe a small proportion of hops longer than the chosen truncation cut-off. When the cut-off value is close to the Lévy stable σ parameter, as it is in our anomalous diffusion models, we observed that the tails of the truncated distribution tend to be undersampled. In order to maintain the distribution's approximate shape up to the cut-off value we recommend ensuring that the cut-off value is at least 2 times σ . However, this may lead to a slight over-prediction of the MSD.

S4.2 Achieving the right correlation structure

We simulated FLM using the algorithm of Stoev and Taqqu [1]. There are no known exact methods for simulating FLM. As a consequence, passing a value of H and α to the algorithm does not necessarily result in the correct correlation structure, although the marginal Lévy stable distribution is correct. We applied a database-based empirical correction in order to use the algorithm to achieve the correct marginal distribution and correlation structure.

Stoev and Taqqu note that the transition between negatively and positively correlated draws occurs when $H = 1/\alpha$. When $\alpha = 2$, the marginal distribution is Gaussian and the transition occurs at $H = 0.5$ as expected from FBM. We corrected the input H so that the value of H measured based on the output sequence equaled the desired H . We first adjusted the value of H by adding $(1/\alpha - 0.5)$, effectively recentering the correlation sign transition for any value of $1 \leq \alpha \leq 2$. This correction alone does a good job for input H values near 0.5, but is insufficient if one desires a low value of H . The exact correction to H is not obvious so we created a database of output H values tabulated as a function of input H and α values. Figure S5 demonstrates the results of applying our correction. Without the correction, FLM realizations are more negatively correlated. This would result in under-predicted mean squared displacements when applying the model.

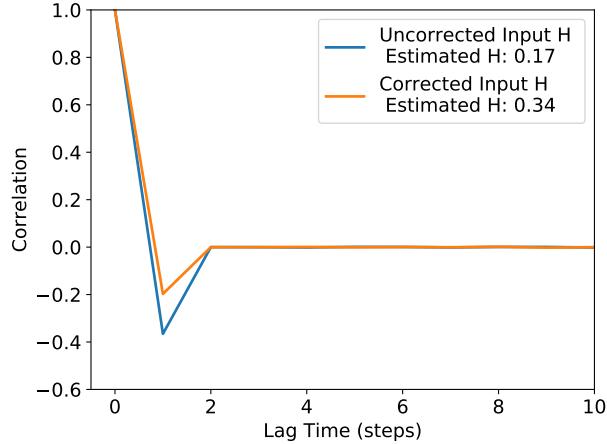


Figure S5: Correcting the Hurst parameter input to the algorithm of Stoev and Taqqu results in an FLM process with a more accurate correlation structure. We generated sequences with an input H of 0.35. We estimated H by fitting the autocorrelation function. Without the correction, H is underestimated, meaning realizations are more negatively correlated than they should be.

S5 Verifying Markovianity

We verified the Markovianity of our transition matrix, T , in two ways. First we ensured that the process satisfied detailed balance:

$$T_{i,j}P_i(t = \infty) = T_{j,i}P_j(t = \infty) \quad (2)$$

where P is the equilibrium distribution of states. This implies that the number of transitions from state i to j and from state j to i should be equal. Graphical representations of the count matrices show that this is true in Figure S6.

Second, we ensured that the transition matrix did not change on coarser time scales. In Figures S6 and S7, we show that increasing the length of time between samples does not change the properties of the count or probability transition matrices.

								Timestep = 0.5 ns								Timestep = 1.0 ns								Timestep = 2.0 ns							
1	22894	485	10529	317	2534	60	1022	45	11424	238	5332	170	1289	31	528	24	5708	121	2720	85	654	16	264	12							
2	500	1030	286	639	63	72	51	54	259	510	136	303	28	39	26	27	132	241	69	151	15	17	13	11							
3	10503	304	43899	1590	1007	22	3008	175	5217	161	21922	797	527	11	1514	87	2593	86	10947	414	253	5	755	38							
4	303	634	1611	4213	44	51	152	330	127	334	822	2114	21	19	69	151	55	166	404	1055	9	11	34	82							
5	2515	63	1042	40	23358	823	9509	379	1261	34	525	21	11629	396	4846	181	652	18	270	10	5784	195	2441	89							
6	59	81	25	54	799	1166	383	485	36	36	13	31	401	584	194	261	15	20	5	15	211	287	97	123							
7	1064	39	2947	168	9541	374	24362	1274	510	23	1453	90	4744	173	12167	604	245	10	722	54	2340	76	6137	294							
8	43	60	172	318	384	483	1281	1938	20	32	77	156	201	243	666	975	10	18	45	79	113	121	333	467							
	1	2	3	4	5	6	7	8	1	2	3	4	5	6	7	8 <th>1</th> <td>2</td> <td>3</td> <td>4</td> <td>5</td> <td>6</td> <td>7</td> <td>8</td>	1	2	3	4	5	6	7	8							

(a) Urea

								Timestep = 0.5 ns								Timestep = 1.0 ns								Timestep = 2.0 ns							
1	9242	7639	1311	888	1709	1540	236	156	4684	3774	666	432	869	742	127	77	2299	1921	354	219	407	377	66	37							
2	7609	25746	958	3955	1363	4266	175	553	3802	12854	493	1997	691	2175	86	311	1907	6428	246	1032	345	1093	43	147							
3	1292	990	886	583	197	122	91	65	630	501	435	280	102	58	50	30	310	259	213	132	58	27	25	21							
4	949	3909	582	3424	118	509	52	458	466	1947	314	1708	57	283	24	237	228	940	148	860	27	148	13	115							
5	1712	1362	201	110	14763	10336	1727	1061	833	682	92	48	7386	5129	887	532	402	352	47	29	3583	2550	470	258							
6	1537	4219	125	533	10396	36209	1083	4598	760	2100	65	261	5204	18121	538	2332	355	1098	31	134	2673	9080	259	1174							
7	239	163	102	58	1705	1103	693	415	116	81	46	22	858	537	347	204	61	44	25	9	434	262	163	100							
8	142	599	62	450	1015	4617	421	3007	59	277	29	217	516	2274	208	1505	34	135	14	113	263	1123	114	750							
	1	2	3	4	5	6	7	8	1	2	3	4	5	6	7	8 <th>1</th> <td>2</td> <td>3</td> <td>4</td> <td>5</td> <td>6</td> <td>7</td> <td>8</td>	1	2	3	4	5	6	7	8							

(b) Ethylene Glycol

Figure S6: The number of transitions from state i to j and j to i are very close indicating that our process obeys detailed balance. Detailed balance is conserved for different sized time steps.

Timestep = 0.5 ns								Timestep = 1.0 ns								Timestep = 2.0 ns							
1	2	3	4	5	6	7	8	1	2	3	4	5	6	7	8	1	2	3	4	5	6	7	8
0.60	0.01	0.28	0.01	0.07	0.00	0.03	0.00	0.60	0.01	0.28	0.01	0.07	0.00	0.03	0.00	0.60	0.01	0.28	0.01	0.07	0.00	0.03	0.00
0.19	0.38	0.11	0.24	0.02	0.03	0.02	0.02	0.20	0.38	0.10	0.23	0.02	0.03	0.02	0.02	0.20	0.37	0.11	0.23	0.02	0.03	0.02	0.02
0.17	0.01	0.73	0.03	0.02	0.00	0.05	0.00	0.17	0.01	0.73	0.03	0.02	0.00	0.05	0.00	0.17	0.01	0.73	0.03	0.02	0.00	0.05	0.00
0.04	0.09	0.22	0.57	0.01	0.01	0.02	0.04	0.03	0.09	0.22	0.58	0.01	0.01	0.02	0.04	0.03	0.09	0.22	0.58	0.00	0.01	0.02	0.05
0.07	0.00	0.03	0.00	0.62	0.02	0.25	0.01	0.07	0.00	0.03	0.00	0.62	0.02	0.26	0.01	0.07	0.00	0.03	0.00	0.61	0.02	0.26	0.01
0.02	0.03	0.01	0.02	0.26	0.38	0.13	0.16	0.02	0.02	0.01	0.02	0.26	0.38	0.12	0.17	0.02	0.03	0.01	0.02	0.27	0.37	0.13	0.16
0.03	0.00	0.07	0.00	0.24	0.01	0.61	0.03	0.03	0.00	0.07	0.00	0.24	0.01	0.62	0.03	0.02	0.00	0.07	0.01	0.24	0.01	0.62	0.03
0.01	0.01	0.04	0.07	0.08	0.10	0.27	0.41	0.01	0.01	0.03	0.07	0.08	0.10	0.28	0.41	0.01	0.02	0.04	0.07	0.10	0.10	0.28	0.39

(a) Urea

Timestep = 0.5 ns								Timestep = 1.0 ns								Timestep = 2.0 ns							
1	2	3	4	5	6	7	8	1	2	3	4	5	6	7	8	1	2	3	4	5	6	7	8
0.41	0.34	0.06	0.04	0.08	0.07	0.01	0.01	0.41	0.33	0.06	0.04	0.08	0.07	0.01	0.01	0.40	0.34	0.06	0.04	0.07	0.07	0.01	0.01
0.17	0.58	0.02	0.09	0.03	0.10	0.00	0.01	0.17	0.57	0.02	0.09	0.03	0.10	0.00	0.01	0.17	0.57	0.02	0.09	0.03	0.10	0.00	0.01
0.31	0.23	0.21	0.14	0.05	0.03	0.02	0.02	0.30	0.24	0.21	0.13	0.05	0.03	0.02	0.01	0.30	0.25	0.20	0.13	0.06	0.03	0.02	0.02
0.09	0.39	0.06	0.34	0.01	0.05	0.01	0.05	0.09	0.39	0.06	0.34	0.01	0.06	0.00	0.05	0.09	0.38	0.06	0.35	0.01	0.06	0.01	0.05
0.05	0.04	0.01	0.00	0.47	0.33	0.06	0.03	0.05	0.04	0.01	0.00	0.47	0.33	0.06	0.03	0.05	0.05	0.01	0.00	0.47	0.33	0.06	0.03
0.03	0.07	0.00	0.01	0.18	0.62	0.02	0.08	0.03	0.07	0.00	0.01	0.18	0.62	0.02	0.08	0.02	0.07	0.00	0.01	0.18	0.61	0.02	0.08
0.05	0.04	0.02	0.01	0.38	0.25	0.15	0.09	0.05	0.04	0.02	0.01	0.39	0.24	0.16	0.09	0.06	0.04	0.02	0.01	0.40	0.24	0.15	0.09
0.01	0.06	0.01	0.04	0.10	0.45	0.04	0.29	0.01	0.05	0.01	0.04	0.10	0.45	0.04	0.30	0.01	0.05	0.01	0.04	0.10	0.44	0.04	0.29

(b) Ethylene Glycol

Figure S7: As the timestep between observations increases, the probability transition matrix does not change significantly.

S6 Derivation of Passage Time Distributions

To derive an analytical equation for the mean first passage time (Equation 11 of the main text), first consider an initial pulse spreading out over time with a fixed mean. We can solve for the time-dependent probability density of particle positions, p , by solving the one dimensional diffusion equation:

$$\frac{\partial p}{\partial t} = D \frac{\partial^2 p}{\partial z^2} \quad (3)$$

The appropriate initial and boundary conditions are:

$$BC1 : t > 0, z = \infty, p = 0$$

$$BC2 : t > 0, z = 0, \frac{\partial p}{\partial z} = 0$$

$$IC : t = 0, c = \delta(z)$$

It has been shown elsewhere that the solution to this equation is: [2]

$$p(z, t) = \frac{1}{\sqrt{4\pi Dt}} \exp\left(\frac{-z^2}{4Dt}\right) \quad (4)$$

We can make the substitution $z = z - vt$, where v represents a constant average velocity, in order to linearly shift the mean as a function of time:

$$p(z, t) = \frac{1}{\sqrt{4\pi Dt}} \exp\left(\frac{-(z - vt)^2}{4Dt}\right) \quad (5)$$

One can track the fraction of particles, F , that have crossed the pore boundary by integrating:

$$F(t) = \int_L^\infty p \, dz = \operatorname{erfc}\left(\frac{L - vt}{2\sqrt{Dt}}\right) \quad (6)$$

where L is the pore length. This represents the cumulative first passage time distribution so we take its derivative in order to arrive at the first passage time distribution:

$$P(t) = -\frac{1}{\sqrt{\pi}} e^{-(L-vt)^2/(4Dt)} \left(-\frac{D(L-vt)}{4(Dt)^{3/2}} - \frac{v}{2\sqrt{Dt}} \right) \quad (7)$$

where the only free parameters for fitting are v and D . We calculated the expected value of Equation 7 in order to get the MFPT. Specifically, we used the python package `scipy.integrate.quad` to numerically integrate:

$$E[t] = \int_0^\infty t P(t) dt \quad (8)$$

S7 Solute hopping and trapping behavior

Analogous to Figure 3 of the main text, Figure S8 demonstrates that all solutes exhibit the same kind of anti-correlated hopping and trapping behavior.

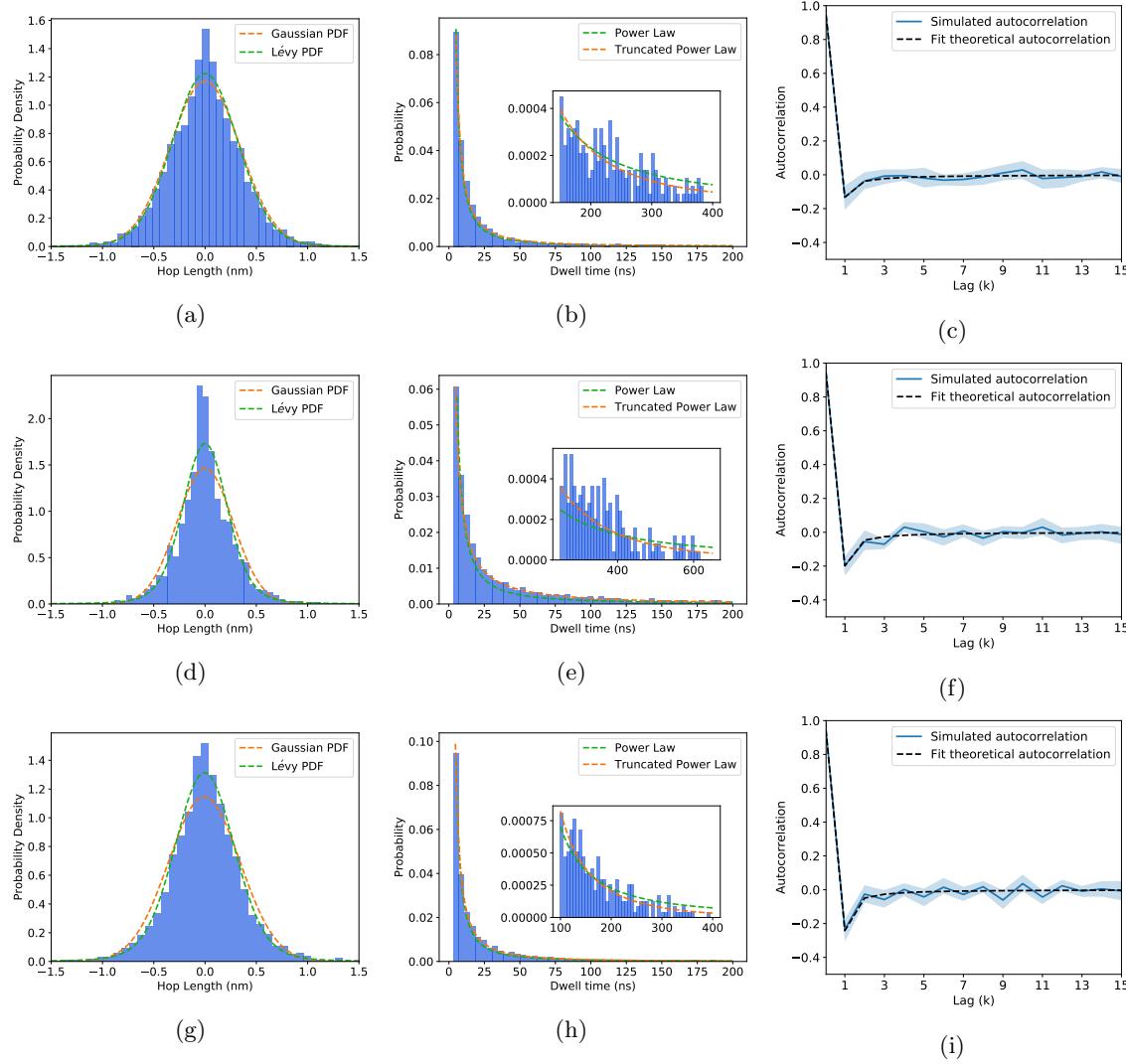


Figure S8: Hop length distributions, dwell time distributions and hop autocorrelation functions respectively for ethylene glycol (a-c), acetic acid (d-f) and methanol (g-i). See Section 3.1.1 of the main text for a more detailed discussion.

S8 AD model MSD Predictions with Pure Power Law Dwell Times

When we use a pure power law distribution to parameterize the dwell time distributions of the 1 and 2 mode AD models, the MD MSDs are severely under-predicted because we are incorporating dwell times on the order of the simulation length into simulated trajectories (see Figure S9). The parameters of the pure power law distribution are included in Figure S10.

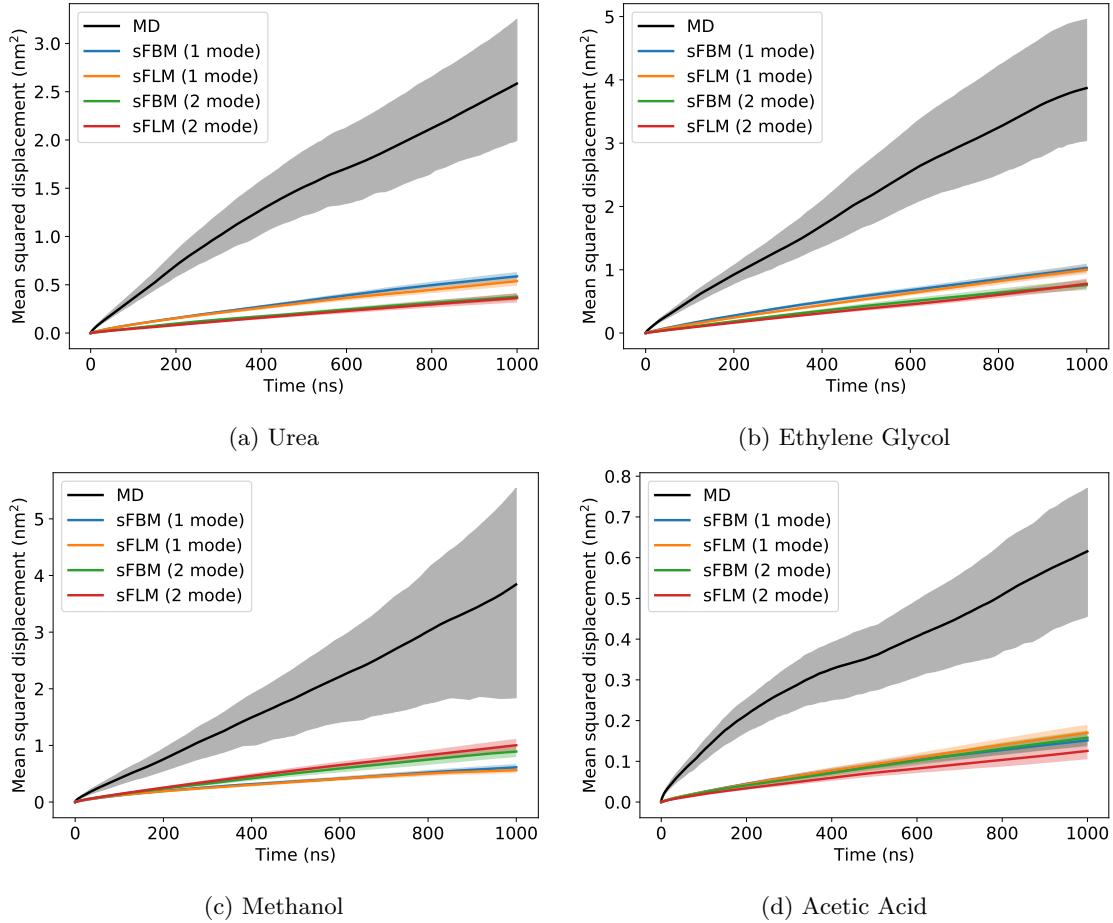


Figure S9: When we do not apply an exponential cut-off to the power law distribution of dwell times, MSDs are consistently under-predicted by the AD model.

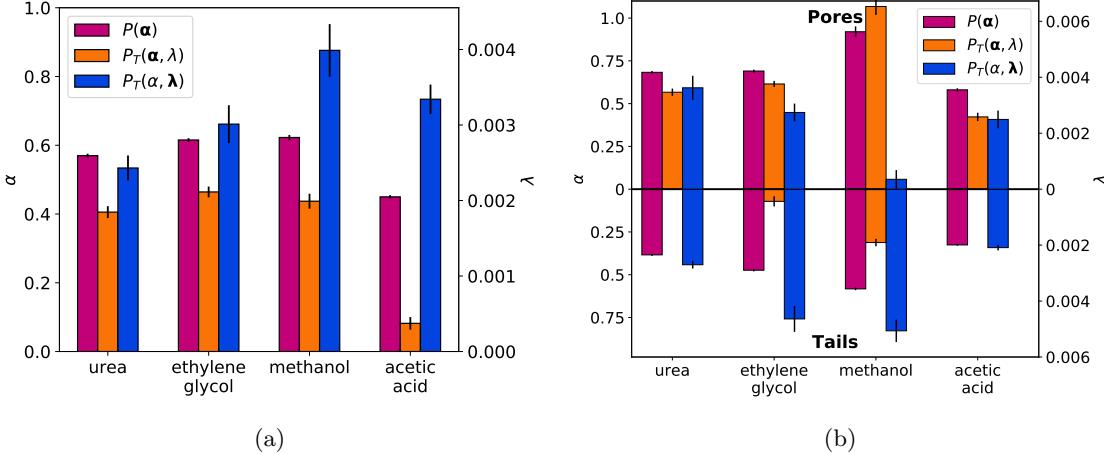


Figure S10: We can parameterize the dwell time distribution in two ways: as a pure power law ($P(\alpha)$) and as a power law with an exponential cut-off ($P(\alpha, \lambda)$). Pure power laws have an infinite variance which allows extremely long dwell times to be sampled (see Figure S9).

S9 Stationarity of Solute Trajectories

We evaluated the predictive capabilities of the AD modeling approach by training the model parameters on the first half of the equilibrated MD trajectory data and then comparing the MSD calculated from AD model realizations to the MSD calculated from the second half of the equilibrated MD trajectory data. This metric is only meaningful if the ensemble of solute trajectories is stationary. In Figure S11, we show that urea and acetic acid show acceptable stationary behavior while methanol and ethylene glycol do not.

We validated both the 1 and 2 mode AD models with urea and acetic acid, since their trajectories appear stationary. The MSDs resulting from 1000 realizations of the AD model are shown in Figure S12. We consider the model's prediction to match well if the MSD lies within the 1σ confidence intervals of the MD MSDs. We also look for qualitative agreement in the shape of the curves.

The models are capable of reasonably predicting the MD MSD values of the second half of the solute trajectories based on parameters generated from the first half when the dwell time distributions are parameterized by a power law with an exponential cut-off. At long timescales, the MSD of Urea is under-predicted for both the 1 and 2 mode models with the same true of acetic acid on short timescales. Without truncation of the power law distribution, the MD MSDs are underestimated in all cases because dwell times on the order of the MD simulation length are sampled and incorporated into the simulated anomalous diffusion trajectories. Considering the longest observed dwell time among all solutes was $1.3 \mu\text{s}$ by ethylene glycol, we believe truncation is well-justified.

From a qualitative perspective, the models only have moderate success at predicting the shape of the MSD curves. Error in the MSD curvature can also help explain some of the error in the predicted MSD magnitudes. The curves predicted for Urea with both 1 and 2 modes appear to have too much curvature, which causes it to under-predict the MSD at long timescales, while those of acetic acid lack curvature, leading the AD model to under-predict the MSD at short timescales.

The under-estimate of Urea's MSD at long timescales is due to long timescale positional anti-correlation which may not be present in the molecularly detailed simulations of the system. The persistent curvature of the predicted MSD curves is a direct consequence of the Hurst parameter. Without anti-correlation, the process would be a pure CTRW for which one would expect the time averaged MSD curves to become linear. [3] It may be true that on the μs time scale, positional correlation is lost which would manifest as a transition from a sub-linear to linear MSD curve. A solution for more accurately modeling this behavior in the future may be to truncate the positional autocorrelation function. [4]

The under-estimate of the curvature of acetic acid's predicted MSD suggests that, in this case, the AD model over-estimates the Hurst parameter. This is not surprising because the Hurst parameter is challenging to quantify, especially with a relatively small amount of data (see Section 2.2.1 of the main text for further

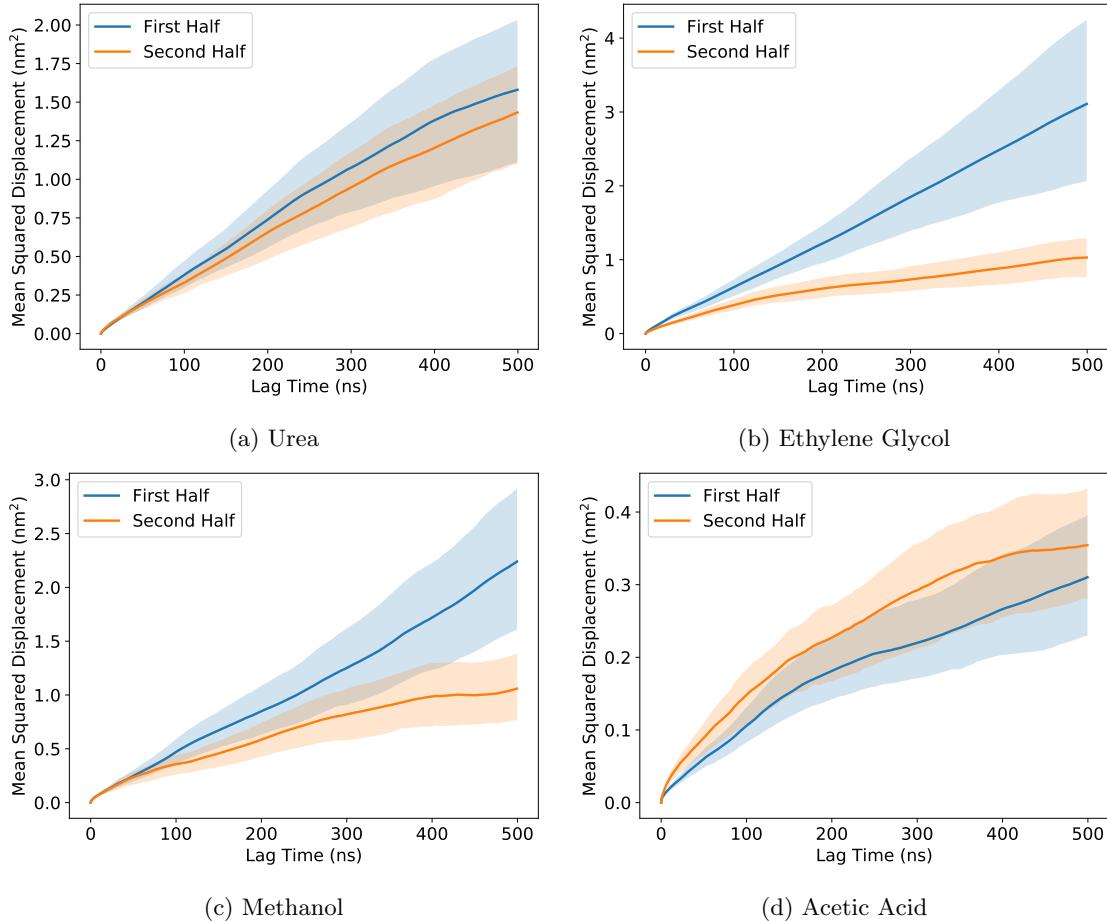


Figure S11: The ensemble of solute trajectories may be stationary if the MSD calculated from different portions of the trajectory are the same. Here we plot the MSD calculated up to a 500 ns time lag of the first and second halves of the equilibrated solute trajectories. Urea and acetic acid have similar MSDs, providing evidence of stationarity, while the MSDs of ethylene glycol and methanol are different suggesting that they are not.

discussion of this challenge). A more accurate measurement of H would fix the shape of the MSD curve, but also lower the predicted MSD meaning we are either underestimating the width of the hop length distribution, favoring longer dwell times, or both.

This brief qualitative analysis suggests two shortcomings of the AD model. First, in a real system, positional anti-correlation may dissipate after a sufficiently long time lag, dependent on the solute being studied. Second, it is difficult to reliably parameterize the Hurst parameter which is important for accurately describing the curvature of the solute MSDs.

This analysis also suggests that working with only half of the data we collected ($\sim 2 \mu\text{s}$ post-equilibration) is not always sufficient for extracting reliable parameter estimates. In most cases, the magnitude of the MSD predictions after a 500 ns time lag are within or close to within error of the MD MSDs (for sFBMcut and sFLMcut), but still appear to systematically under-predict the mean. We may be operating on the border of the minimum amount of data required to accurately parameterize the AD model. Longer simulations and more independent trajectories may be necessary. Therefore, in the next section we will work with parameters fit to the full equilibrated portion of the solute trajectories, doubling the data.

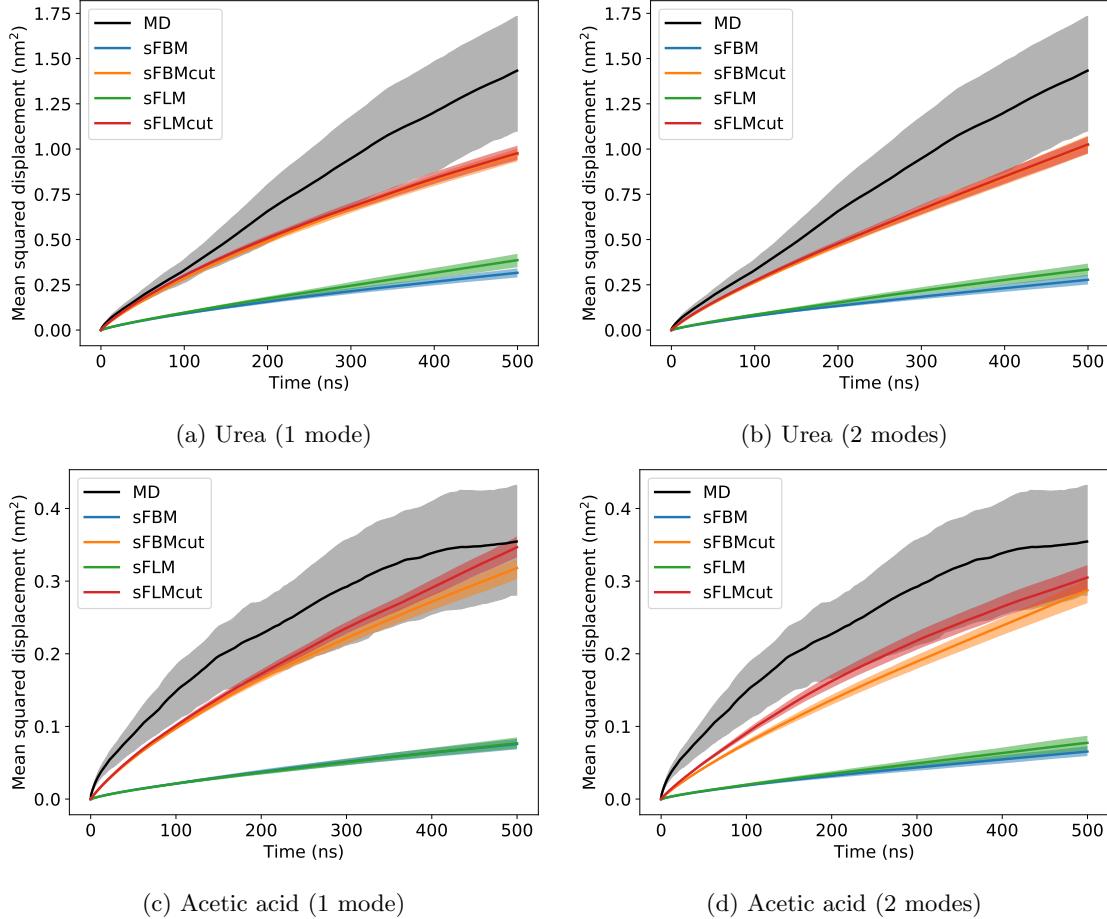


Figure S12: In most cases, when using power laws with exponential cut-offs (sFBMcut and sFLMcut), the MSD curves predicted by the AD model trained on the first half of the equilibrated data lie within the 1σ confidence intervals of the MD MSD curves generated from the second half of the equilibrated solute data. The models which use pure power laws systematically under-predict the MD MSD curves. Over and under-estimated curvature of the of urea and acetic acid's MSD curves respectively causes the magnitude of urea's predicted MSDs to be under-predicted at long timescales and those of acetic acid to be under-predicted at short timescales.

S10 Tabular Anomalous Diffusion Parameters

The tables in the section are tabular representations of the parameters depicted in Figures 6 and 7 of the main text.

1 Mode Model	Parameters	Urea	Ethylene Glycol	Methanol	Acetic Acid
Dwell Distributions	$P(\alpha_d)$	0.57	0.62	0.62	0.45
	$P_T(\alpha_d, \lambda)$	0.40, 0.0024	0.47, 0.0030	0.44, 0.0040	0.08, 0.0033
Hop Distributions	$\mathcal{N}(\sigma)$	0.33	0.34	0.35	0.27
	$L(\sigma, \alpha_h)$	0.21, 1.84	0.23, 1.92	0.22, 1.80	0.16, 1.72
Correlation	$\gamma(H)$	0.37	0.40	0.30	0.34

Table S2: To create a 1 mode model for each solute, we parameterized a pure power law ($P(\alpha_d)$) and a truncated power law ($P_T(\alpha_d, \lambda)$) distribution to describe solute dwell times. Lower values of α_d lead to heavier power law tails and higher values of λ truncate the distribution at lower dwell times. We parameterized Gaussian ($\mathcal{N}(\sigma)$) and Lévy stable ($L(\sigma, \alpha_h)$) distributions to describe solute hop lengths. We assume the mean (μ) to be zero for these distributions and there to be no skewness ($\beta = 0$) in the Lévy stable distributions. High values of σ and lower values of α_h result in larger hops. Finally, we parameterized the hop autocorrelation function ($\gamma(H)$) to describe the degree of correlation between hops. Higher values of H display closer to Brownian behavior.

2 Mode Model						
	Parameters	Mode	Urea	Ethylene Glycol	Methanol	Acetic Acid
Dwell Distributions	$P(\alpha_d)$	1	0.69	0.69	0.90	0.58
		2	0.38	0.48	0.58	0.33
	$P_T(\alpha_d, \lambda)$	1	0.56, 0.0037	0.62, 0.0026	1.04, 0.0006	0.41, 0.0026
		2	0.00, 0.0027	0.06, 0.0049	0.30, 0.0054	0.00, 0.0021
Hop Distributions	$\mathcal{N}(\sigma)$	1	0.35	0.38	0.45	0.32
		2	0.24	0.23	0.32	0.17
	$L(\sigma, \alpha_h)$	1	0.24, 1.91	0.26, 1.99	0.31, 1.97	0.21, 1.91
		2	0.12, 1.50	0.15, 1.90	0.20, 1.85	0.09, 1.50
Correlation	$\gamma(H)$	—	0.37	0.40	0.30	0.34

Table S3: The two model parameterizes solute behavior in the pores and tails separately. Generally, movement is much more restricted in the tail region. To create a 2 mode model, we generated a set of parameters based on solute behavior as function of distance from the pore center. Mode 1 corresponds to solute behavior within 0.75 nm of the pore center and mode 2 corresponds to behavior greater than or equal to 0.75 nm from the pore center. Note that we used the same Hurst parameter for both modes due to a low number of sufficiently long sequences of hops in each mode. See Table S2 for descriptions of the parameters.

S11 Tabular MSDDM parameters

The following table is a tabular representation of the parameters depicted in Figure 10 of the main text.

State	Urea			Ethylene Glycol			Methanol			Acetic Acid		
	H	α_h	σ	H	α_h	σ	H	α_h	σ	H	α_h	σ
1	0.10	1.79	0.034	0.09	1.68	0.045	0.11	1.56	0.052	0.10	1.78	0.035
2	0.06	1.80	0.033	0.09	1.75	0.037	0.07	1.63	0.043	0.08	1.88	0.032
3	0.11	1.88	0.030	0.11	1.86	0.030	0.02	1.80	0.036	0.04	2.00	0.030
4	0.10	1.95	0.027	0.04	1.91	0.028	0.02	1.75	0.036	0.04	2.00	0.027
5	0.19	1.34	0.048	0.15	1.40	0.062	0.10	1.28	0.074	0.13	1.47	0.048
6	0.15	1.45	0.040	0.11	1.52	0.040	0.03	1.50	0.042	0.09	1.70	0.038
7	0.15	1.61	0.032	0.05	1.60	0.040	0.28	1.20	0.043	0.08	1.77	0.031
8	0.11	1.71	0.028	0.05	1.74	0.030	0.04	1.83	0.037	0.01	2.00	0.030
T	0.34	1.42	0.036	0.37	1.44	0.045	0.35	1.45	0.057	0.34	1.54	0.040

Table S4: We calculated values of H , α_h and σ from MD simulation trajectories and used them to generate realizations of our MSDDM model. The states are defined in Table 2 of the main text except state T which describes the transition emissions.

Analytical fits to MFPT distributions

In Figures S13 and S14, we demonstrate the high quality of our analytical fits of Equation 7 to the distribution of solute first passage times derived from both the AD and MSDDM models. The histograms in Figure S14 are relatively noisy because we only generated 1000 realizations of the MSDDM versus 10000 of the AD model. In Figure S15, we justify our use of fewer MSDDM realizations by using subsets of the 10000 AD model realizations to show that one needs as few as 100 independent trajectory realizations at each pore length in order to reliably fit Equation 7 to the passage time distributions.

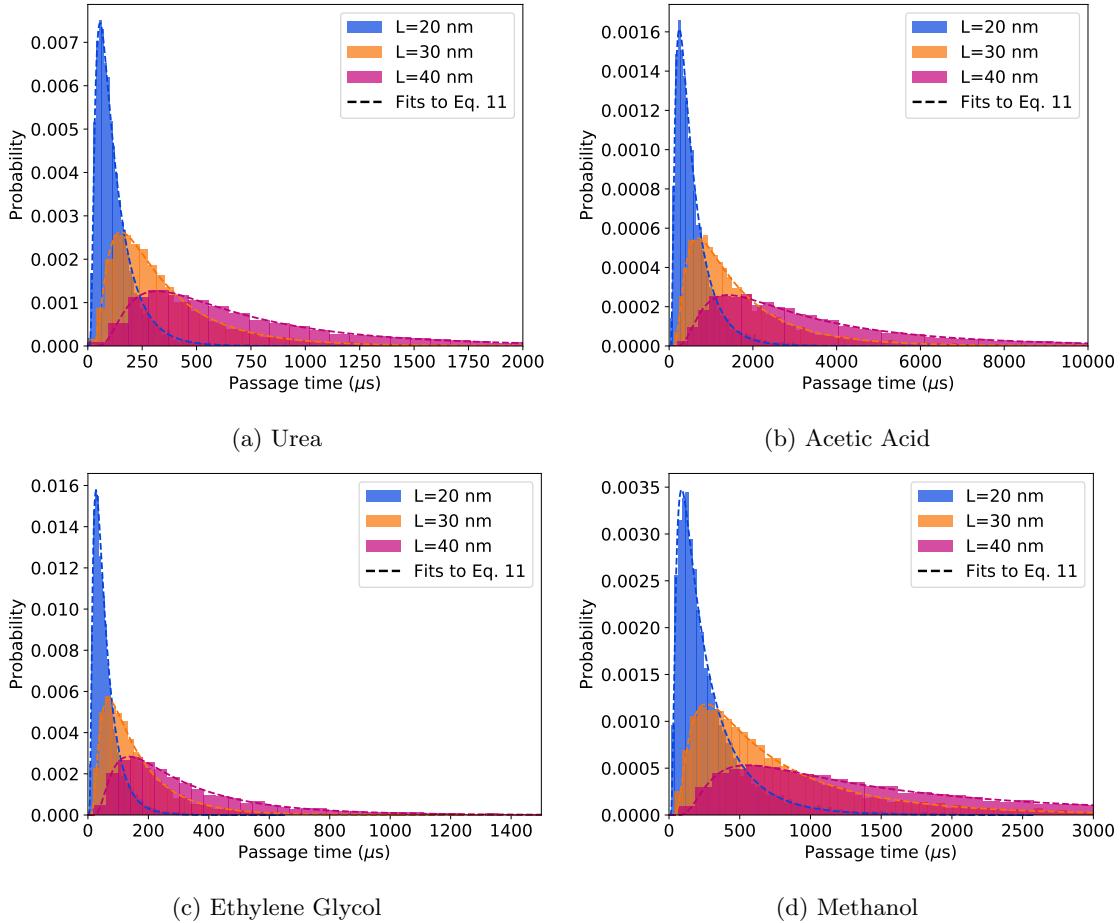


Figure S13: We fit Equation 11 of the main text to the first passage time distributions generated by 10,000 realizations of the anomalous diffusion model.

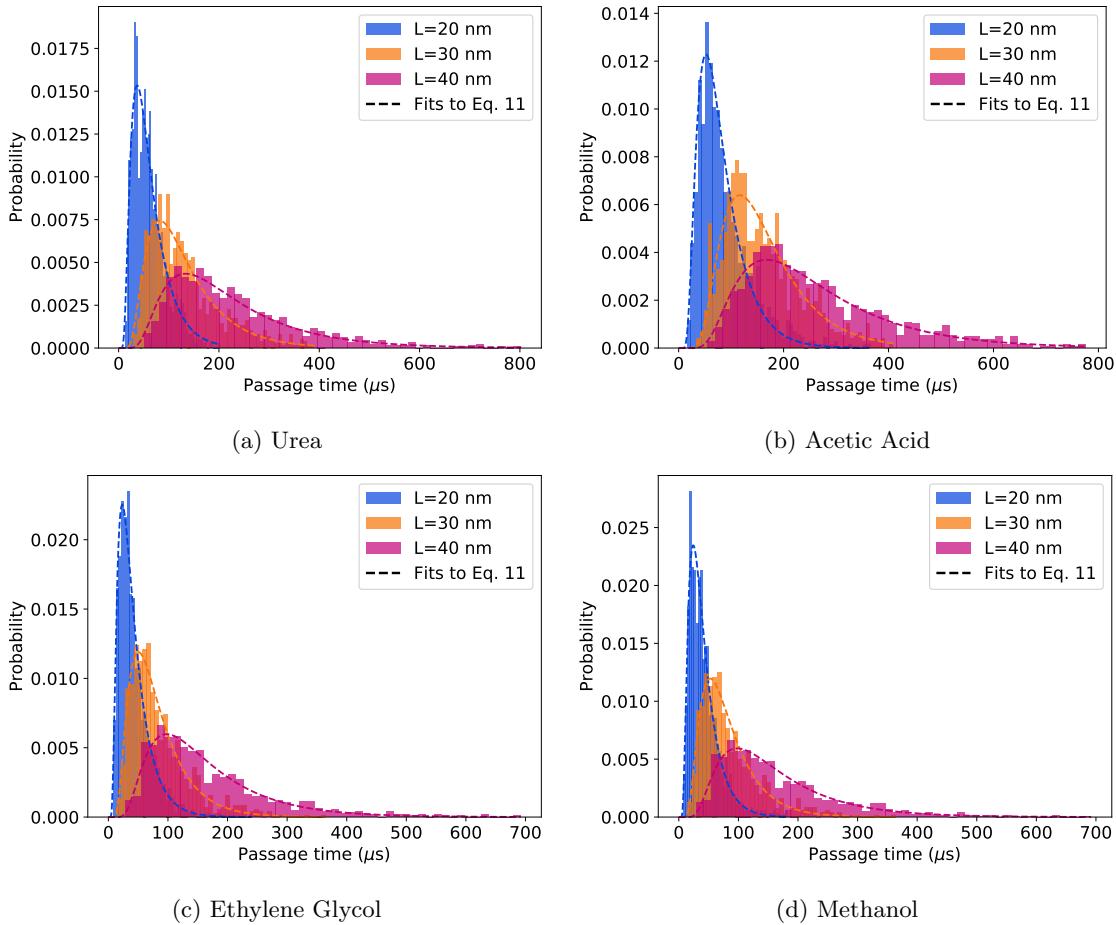


Figure S14: We fit Equation 11 of the main text to the first passage time distributions generated by 1000 realizations of the Markov state dependent dynamical model. Note that we generated 10 times less realizations of the MSDDM which leads to noisier histograms. We show that this has a negligible effect on the fits in Figure S15.

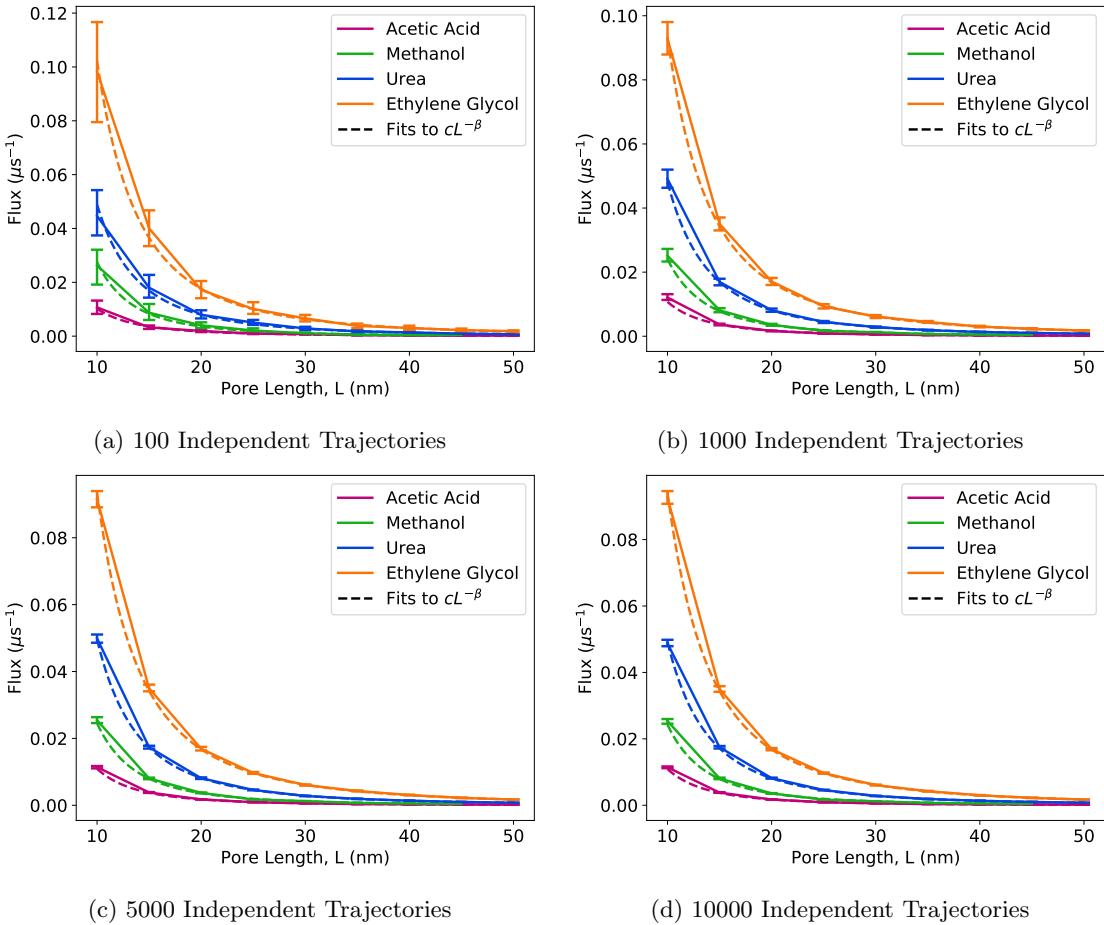


Figure S15: Even using a small number of independent trajectories, one can reliably calculate flux as a function of pore length. The uncertainty in the flux curves decreases as we add more independent trajectories.

S12 Explanation of Under-Estimated MSDDM β Values

The value of β for the MSDDM is under-estimated due to assumptions of the model itself as well as inaccuracies in the correlation structure of very long FLM trajectories. Turning first to the model itself, we have designed it to ensure that the magnitude of the hops in the series of transitions between states are anti-correlated from start to finish. This assumes that the transition correlation structure is unaffected by the sub-trajectories between each state transition. Each time a state transition occurs, one must initialize a new time series sub-trajectory with its own correlation structure. Since we add the transitional hop lengths to each end of trapped state sub-trajectories, the transitional hop lengths are shifted with respect to one another, decreasing correlation between them. We tested this reasoning by modifying the MSDDM to completely immobilize particles except when they transition between states. Surprisingly, β is still close to the Brownian value. Further experimentation reveals that this is actually a consequence of the FLM simulation procedure. Simulating FLM requires Riemann-sum approximations of the stochastic integrals defining the process. To generate the curves in Figure 13c, we needed to correlate 25–1000 times more hops than for the MSD predictions in Figure 11. In short, it is computationally infeasible to use enough terms to accurately incorporate long timescale correlations into our long MSDDM realizations. Thus at long timescales, we lose correlation between transitional jumps. We confirmed this hypothesis by using fractional Brownian motion, for which we have an exact simulation method, in place of FLM in the MSDDM algorithm. When we use FBM, β increases well above the Brownian value. β increases even further if we immobilize the trapped states. Thus the low value of β of the MSDDM is a consequence of model assumptions and inexact simulation of FLM.

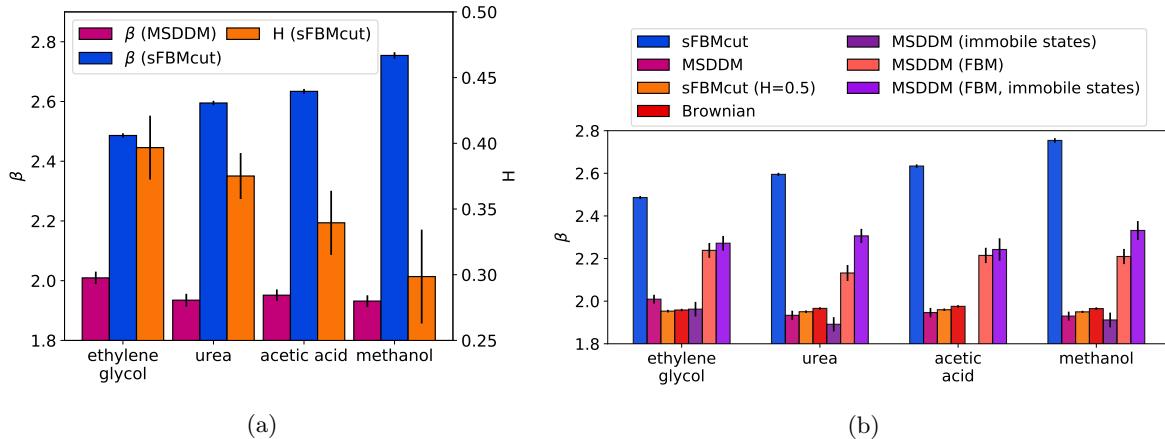


Figure S16: (a) The β values of the sFBMcut model appear to be inversely proportional to the Hurst parameters. The β values of the MSDDM are much lower with much less variation than the sFBMcut model. (b) The β parameters of the MSDDM are low because sub-trajectories between state transitions decorrelate transitional hops and because our FLM simulation procedure does not accurately correlate hops after long time lags. The high β parameters of the sFBMcut model are a consequence of anti-correlation between hops. Removing hop correlations causes β to drop down close to MSDDM values (sFBMcut ($H=0.5$)). Removing dwell times in addition to hop correlation yields a similar value of β (Brownian). Immobilizing particles while in a trapped state also yields a similar value of β (MSDDM (immobile states)). Replacing FLM with FBM in the MSDDM raises β well above the Brownian value (MSDDM (FBM)). Replacing FLM with FBM and immobilizing particles in trapped states further raises the value of β .

References

- [1] S. Stoev and M. S. Taqqu, “Simulation methods for linear fractional stable motion and farima using the fast fourier transform,” *Fractals*, vol. 12, pp. 95–121, Mar. 2004.
- [2] E. L. Cussler, *Diffusion: Mass Transfer in Fluid Systems*. Cambridge University Press, 3 ed., 2009.
- [3] Y. Meroz and I. M. Sokolov, “A Toolbox for Determining Subdiffusive Mechanisms,” *Phys. Rep.*, vol. 573, pp. 1–29, Apr. 2015.
- [4] D. Molina-Garcia, T. Sandev, H. Safdari, G. Pagnini, A. Chechkin, and R. Metzler, “Crossover from anomalous to normal diffusion: truncated power-law noise correlations and applications to dynamics in lipid bilayers,” *New J. Phys.*, vol. 20, p. 103027, Oct. 2018.