

# Uncertainty Quantification for Molecular Property Predictions with Graph Neural Architecture Search

Shengli Jiang<sup>1\*</sup>, Shiyi Qin<sup>1</sup>, Reid C. Van Lehn<sup>1</sup>,  
Prasanna Balaprakash<sup>2</sup> and Victor M. Zavala<sup>1,3</sup>

<sup>1</sup>Department of Chemical and Biological Engineering  
University of Wisconsin-Madison, 1415 Engineering Dr, Madison, WI 53706

<sup>2</sup>Computing and Computational Sciences Directorate  
Oak Ridge National Laboratory, P.O. Box 2008, Oak Ridge, TN 37831

<sup>3</sup>Mathematics and Computer Science Division  
Argonne National Laboratory, Lemont, IL 60439

## S1 AutoGNNUQ

### S1.1 Attention Functions

The attention functions used to calculate  $\alpha_{vw}$  are shown in Table S1. For constant attention,  $\alpha_{vw}$  is always 1. For GAT,  $\alpha_{vw}$  is  $\text{LeakyReLU}(\mathbf{a}(\mathbf{W}\mathbf{h}_v || \mathbf{W}\mathbf{h}_w))$ , where  $\mathbf{a}$  is a trainable vector,  $\mathbf{W}$  is a trainable weight matrix,  $||$  denotes concatenation, and  $\mathbf{h}_v$  and  $\mathbf{h}_w$  are the node hidden feature vectors for nodes  $v$  and  $w$ , respectively. The SYM-GAT attention coefficient is procured by the summation of  $\alpha_{vw}$  and  $\alpha_{wv}$  from GAT. The COS attention type computes  $\alpha_{vw}$  as  $\mathbf{a}(\mathbf{W}\mathbf{h}_v || \mathbf{W}\mathbf{h}_w)$ . In linear attention,  $\alpha_{vw}$  is  $\tanh(\mathbf{a}_l \mathbf{W}\mathbf{h}_v + \mathbf{W}\mathbf{h}_w)$ , with  $\mathbf{a}_l$  being a learnable vector. Generalized linear attention, or gen-linear attention, determines  $\alpha_{vw}$  as  $\mathbf{W}_G \tanh(\mathbf{W}\mathbf{h}_v + \mathbf{W}\mathbf{h}_w)$ , where  $\mathbf{W}_G$  is a trainable matrix.

Table S1: Set of attention functions, where  $||$  denotes the concatenation;  $\mathbf{a}$ ,  $\mathbf{a}_l$ , and  $\mathbf{a}_r$ ; and  $\mathbf{W}_G$  is a trainable matrix.

Attention Mechanism	Equations
Constant	1
GAT	$\text{LeakyReLU}(\mathbf{a}(\mathbf{W}\mathbf{h}_v    \mathbf{W}\mathbf{h}_w))$
SYM-GAT	$\alpha_{vw} + \alpha_{wv}$ based on GAT
COS	$\mathbf{a}(\mathbf{W}\mathbf{h}_v    \mathbf{W}\mathbf{h}_w)$
Linear	$\tanh(\mathbf{a}_l \mathbf{W}\mathbf{h}_v + \mathbf{W}\mathbf{h}_w)$
Gen-Linear	$\mathbf{W}_G \tanh(\mathbf{W}\mathbf{h}_v + \mathbf{W}\mathbf{h}_w)$

\*Corresponding Author: sjiang87@wisc.edu

## S1.2 Architectures Discovered

Figures S1-S4 presents the optimal model architectures identified for the Lipo dataset using a random seed of 0, while Figures S5-S8 depicts the least effective models. In this context, input 0 corresponds to node features, input 1 to edge features, input 2 to edge pairs, and input 3 to node masks. A clear observation is that superior models generally incorporate more skip-connections and predominantly use the global attention pooling function. Conversely, the less effective models exhibit simpler architectures with rudimentary readout functions such as simple flatten.

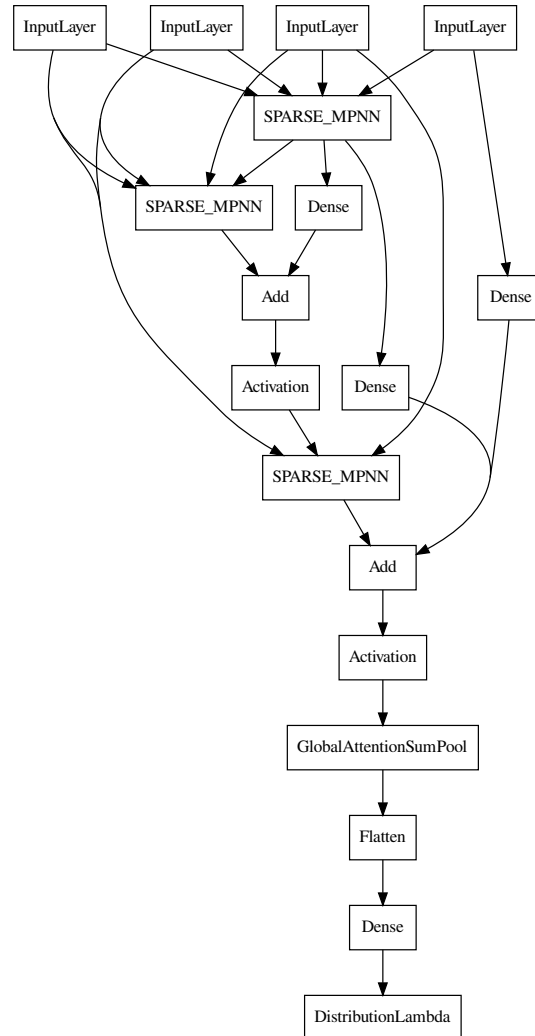


Figure S1: Model architectures demonstrating the highest search rewards (log-likelihood) for the Lipo dataset, using a random seed of 0.

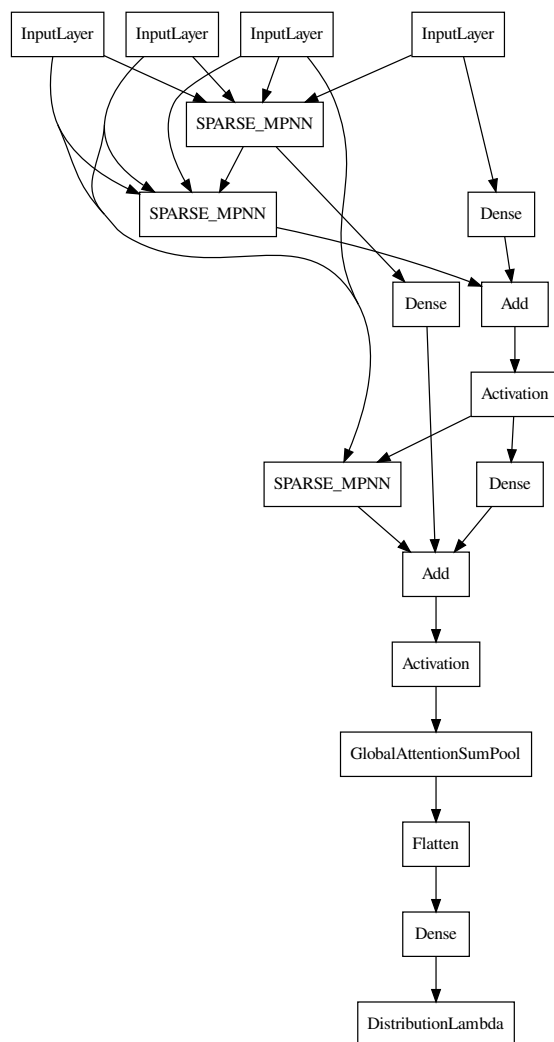


Figure S2: Model architectures demonstrating the second highest search rewards (log-likelihood) for the Lipo dataset, using a random seed of 0.

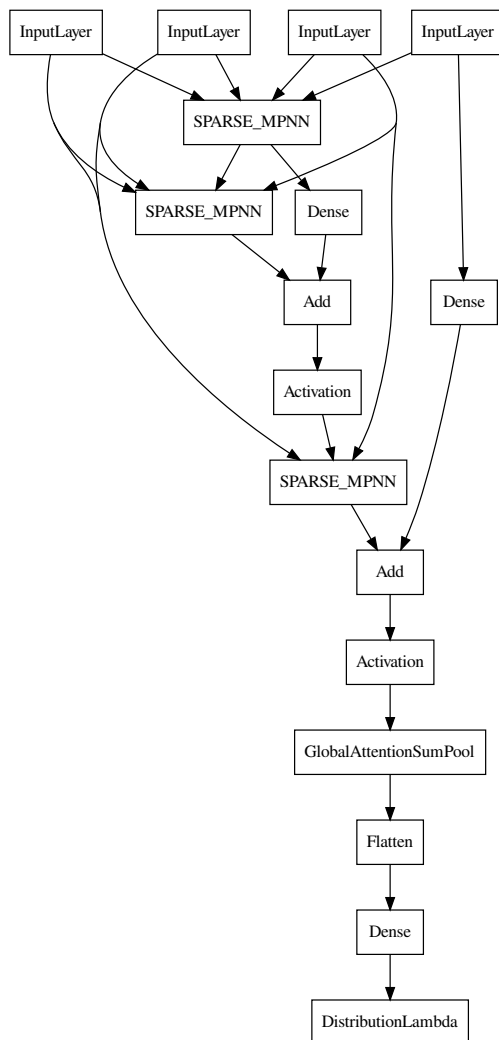


Figure S3: Model architectures demonstrating the third highest search rewards (log-likelihood) for the Lipo dataset, using a random seed of 0.

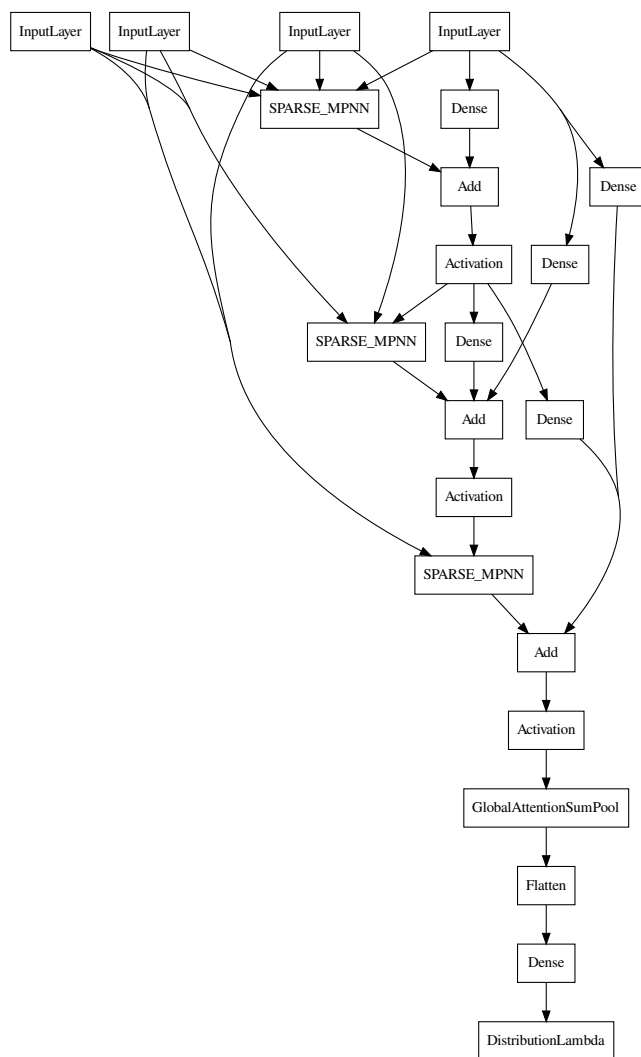


Figure S4: Model architectures demonstrating the fourth highest search rewards (log-likelihood) for the Lipo dataset, using a random seed of 0.

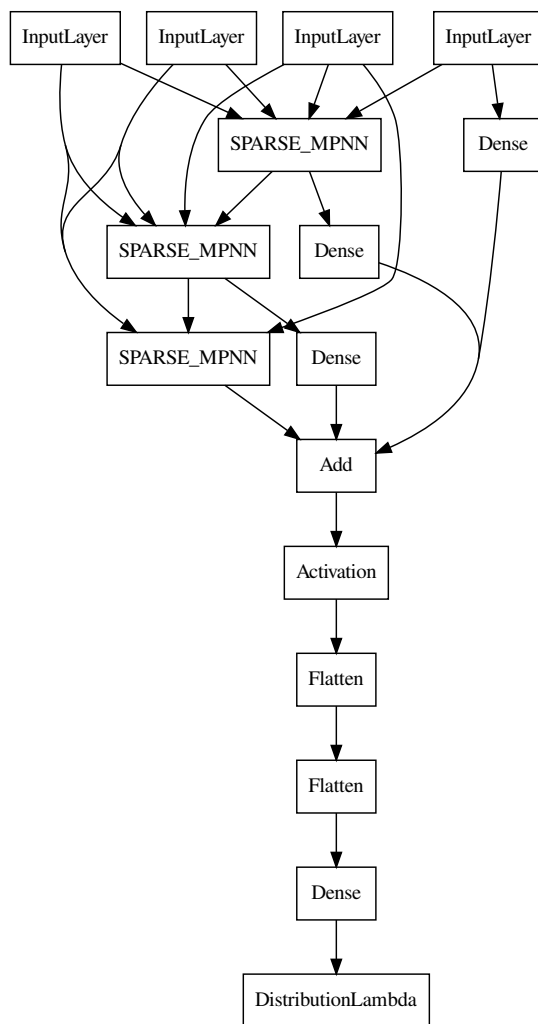


Figure S5: Model architectures demonstrating the lowest search rewards (log-likelihood) for the Lipo dataset, using a random seed of 0.

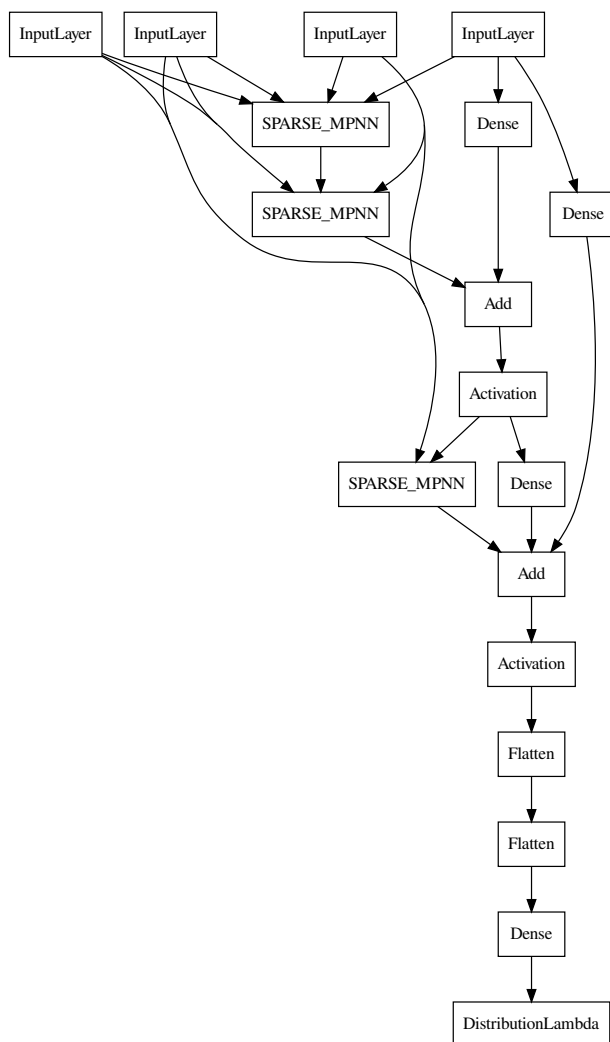


Figure S6: Model architectures demonstrating the second lowest search rewards (log-likelihood) for the Lipo dataset, using a random seed of 0.

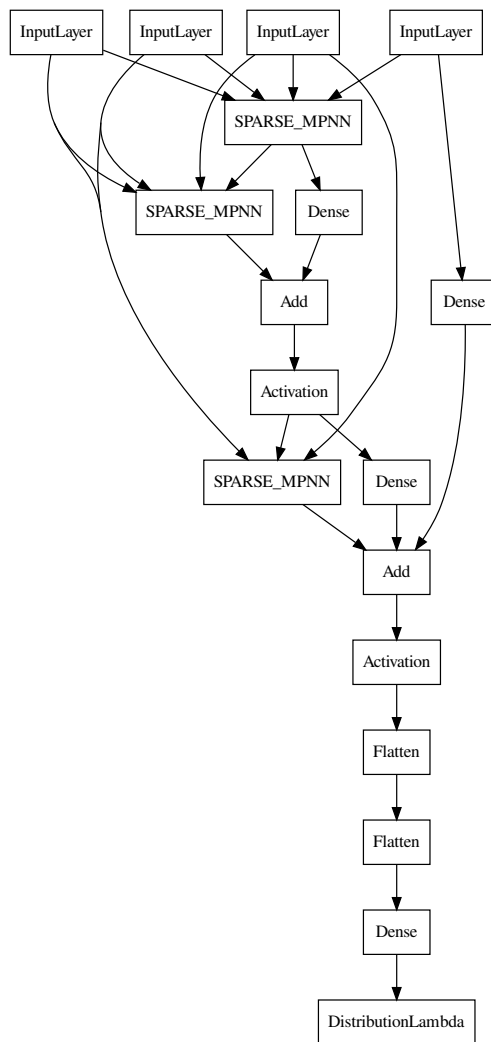


Figure S7: Model architectures demonstrating the third lowest search rewards (log-likelihood) for the Lipo dataset, using a random seed of 0.



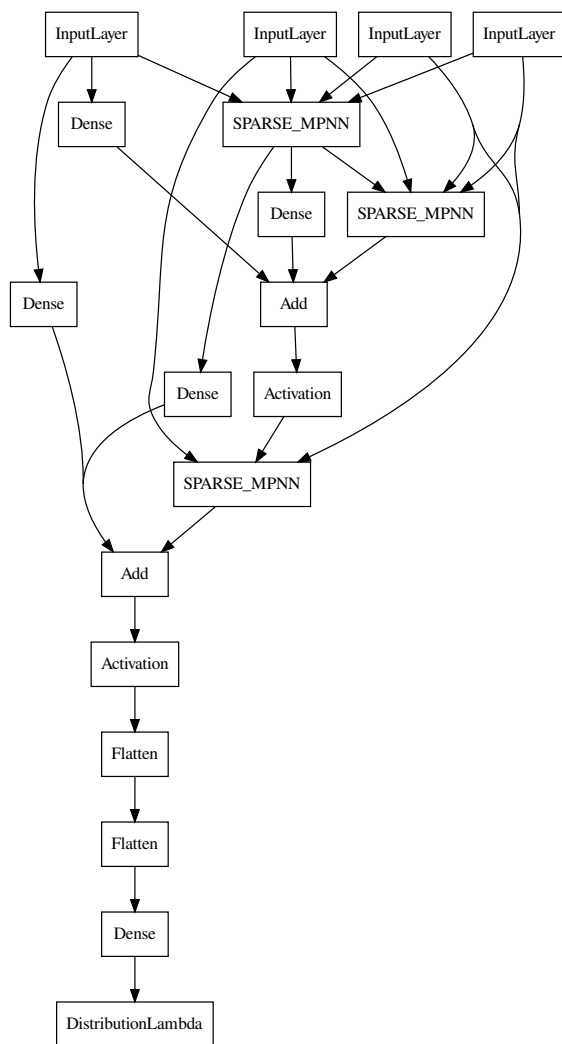


Figure S8: Model architectures demonstrating the fourth lowest search rewards (log-likelihood) for the Lipo dataset, using a random seed of 0.

## S2 Prediction Parity Plots of AutoGNNUQ

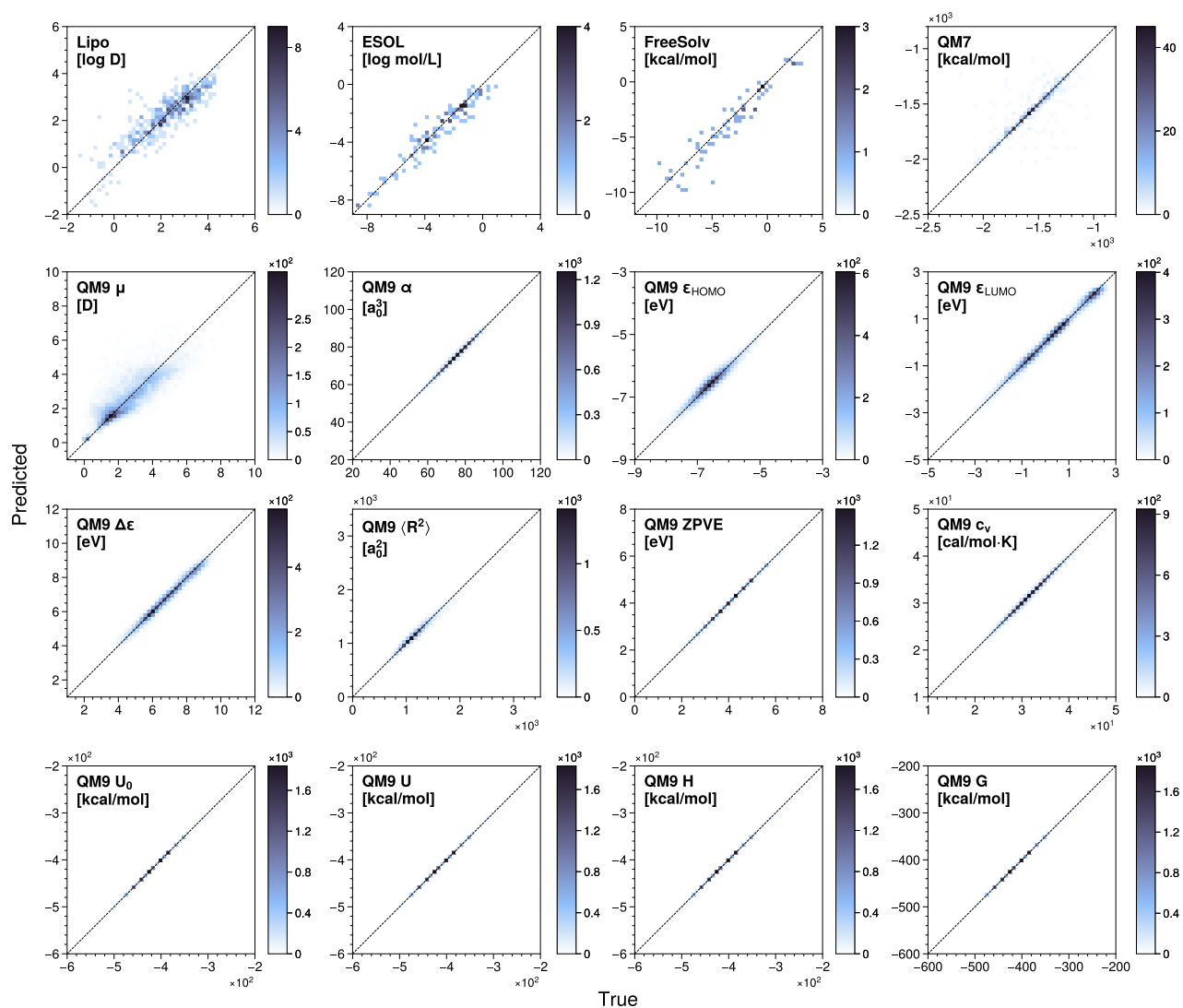


Figure S9: Prediction parity plot with a random seed of 1.

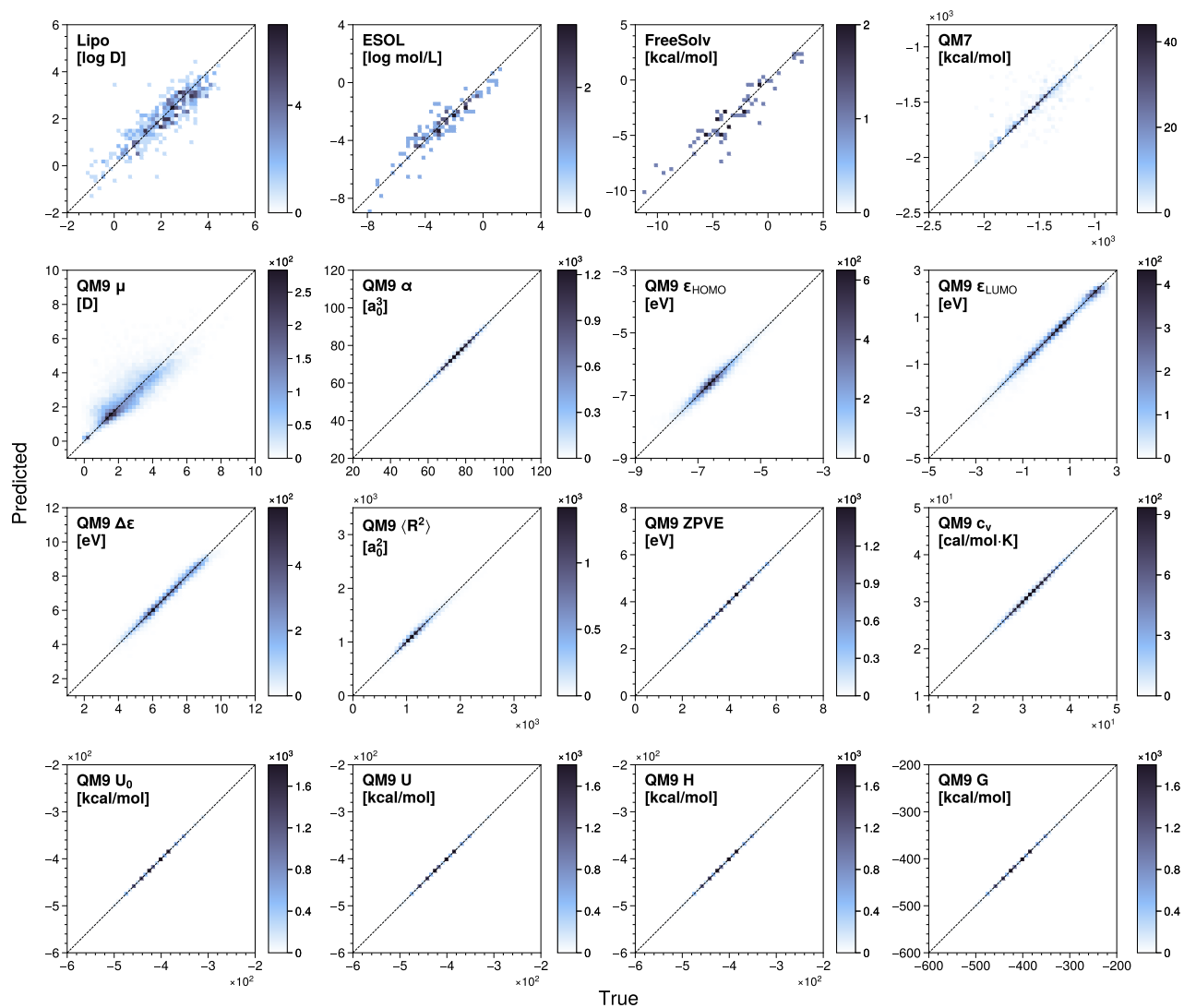


Figure S10: Prediction parity plot with a random seed of 2.

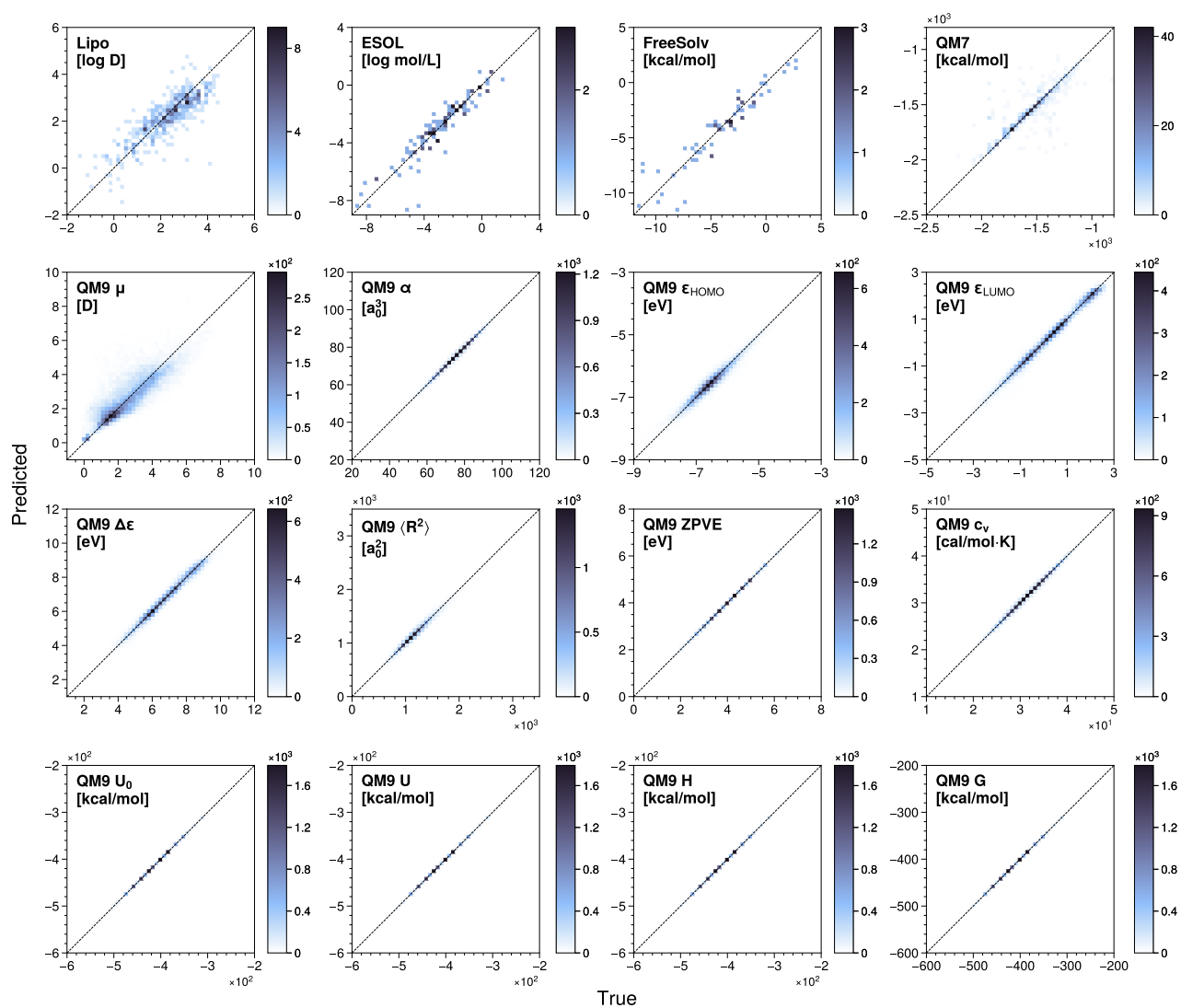


Figure S11: Prediction parity plot with a random seed of 3.

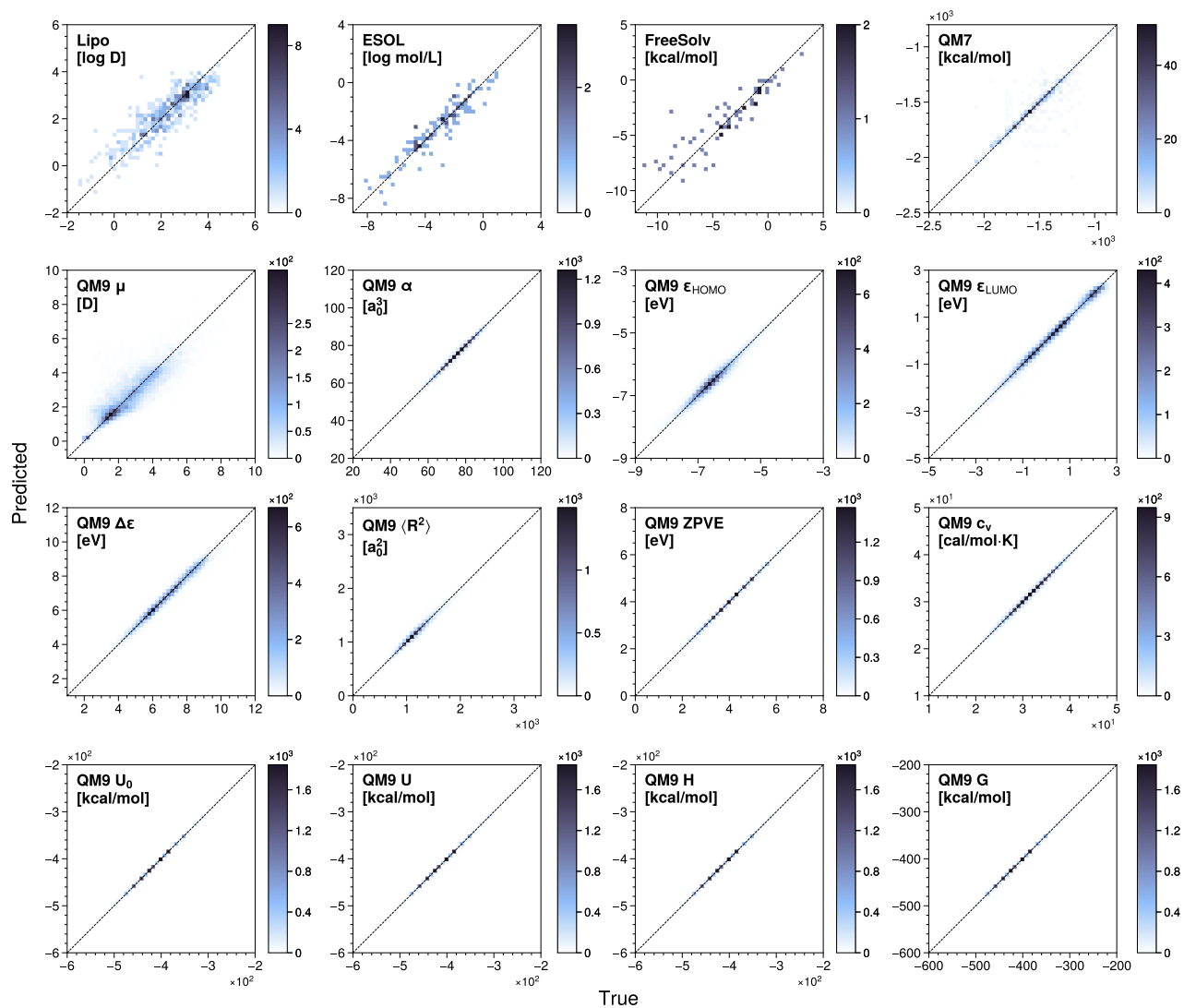


Figure S12: Prediction parity plot with a random seed of 4.

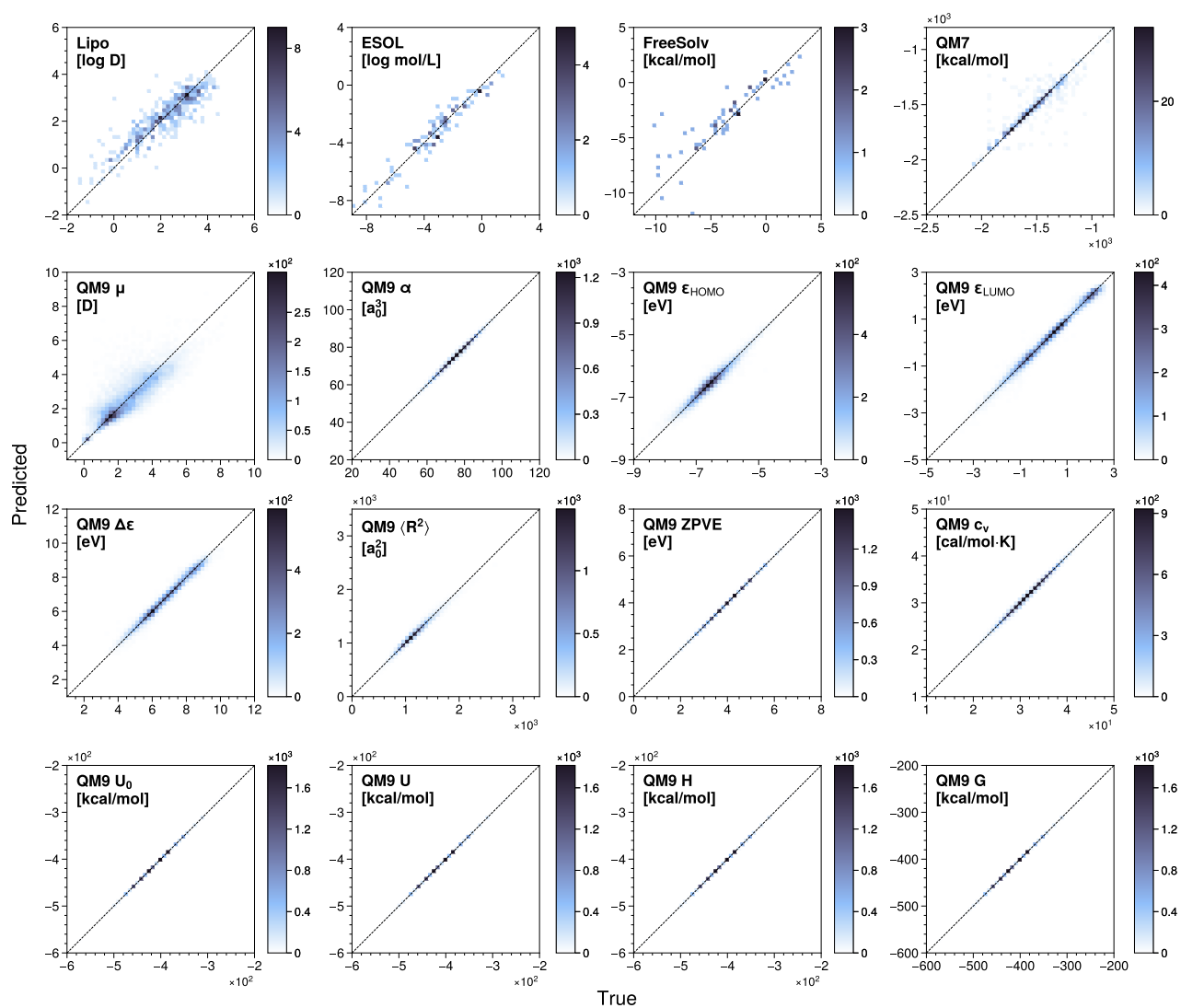


Figure S13: Prediction parity plot with a random seed of 5.

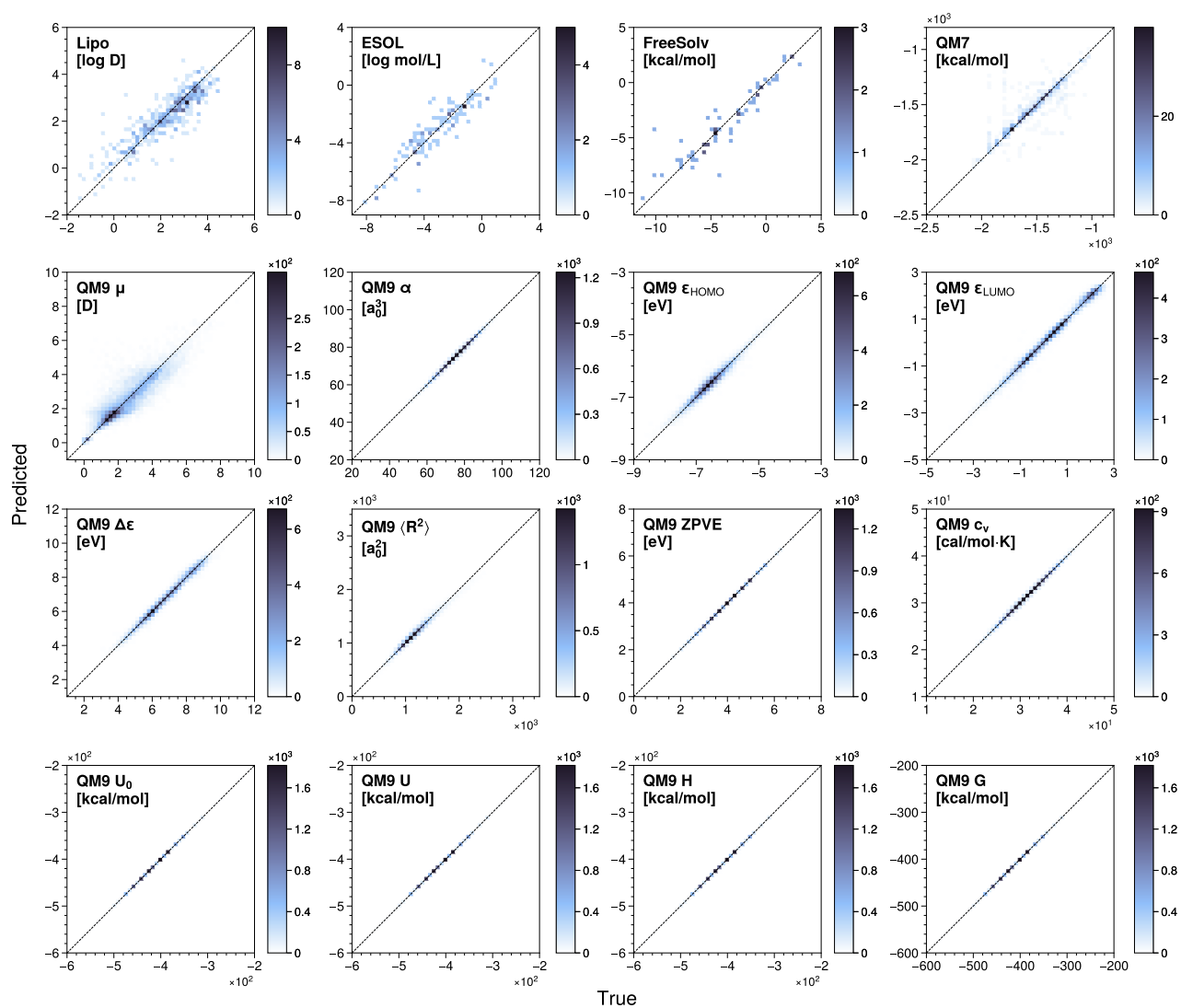


Figure S14: Prediction parity plot with a random seed of 6.

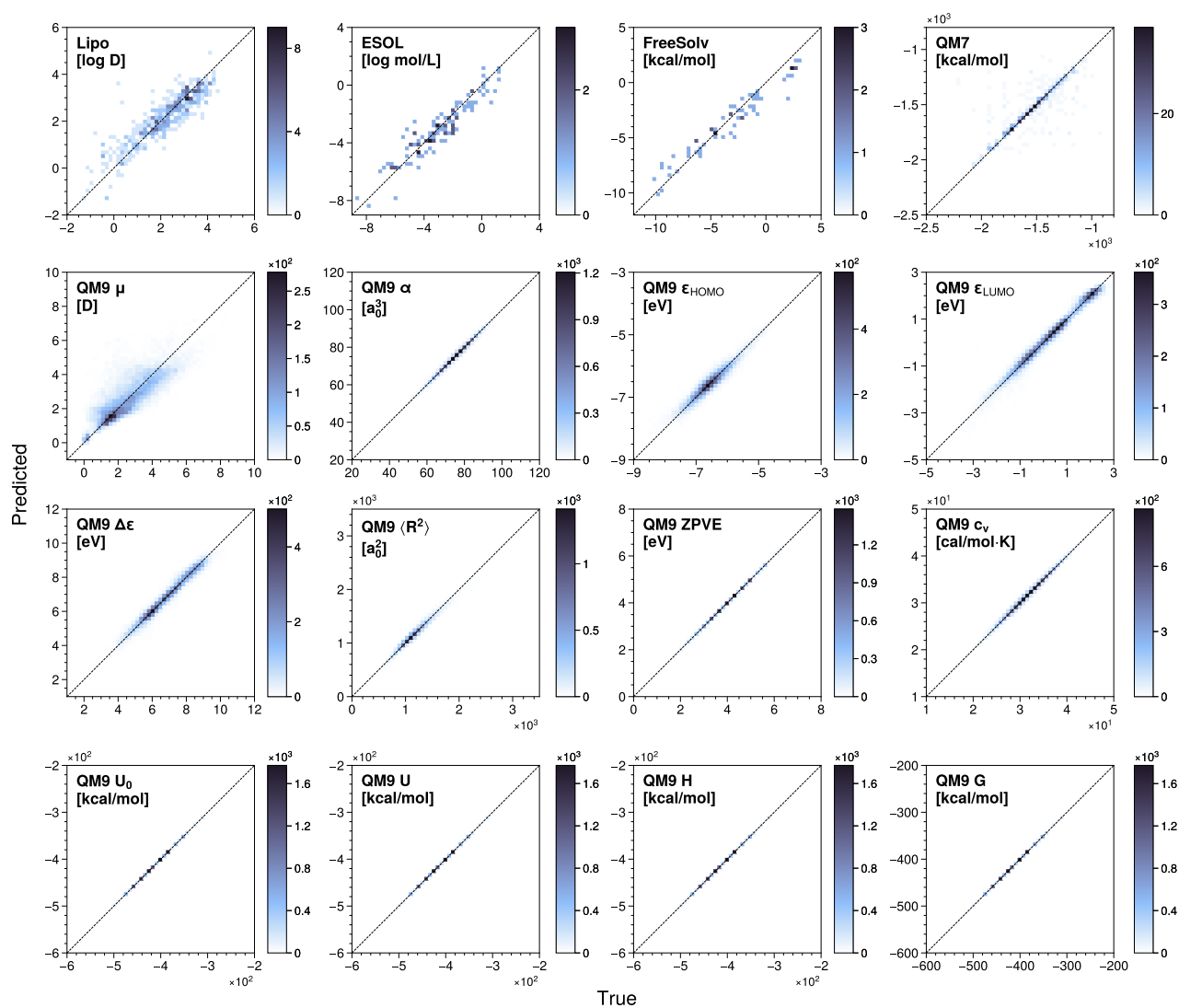


Figure S15: Prediction parity plot with a random seed of 7.



## S3 Prediction Error and Uncertainty of AutoGNNUQ

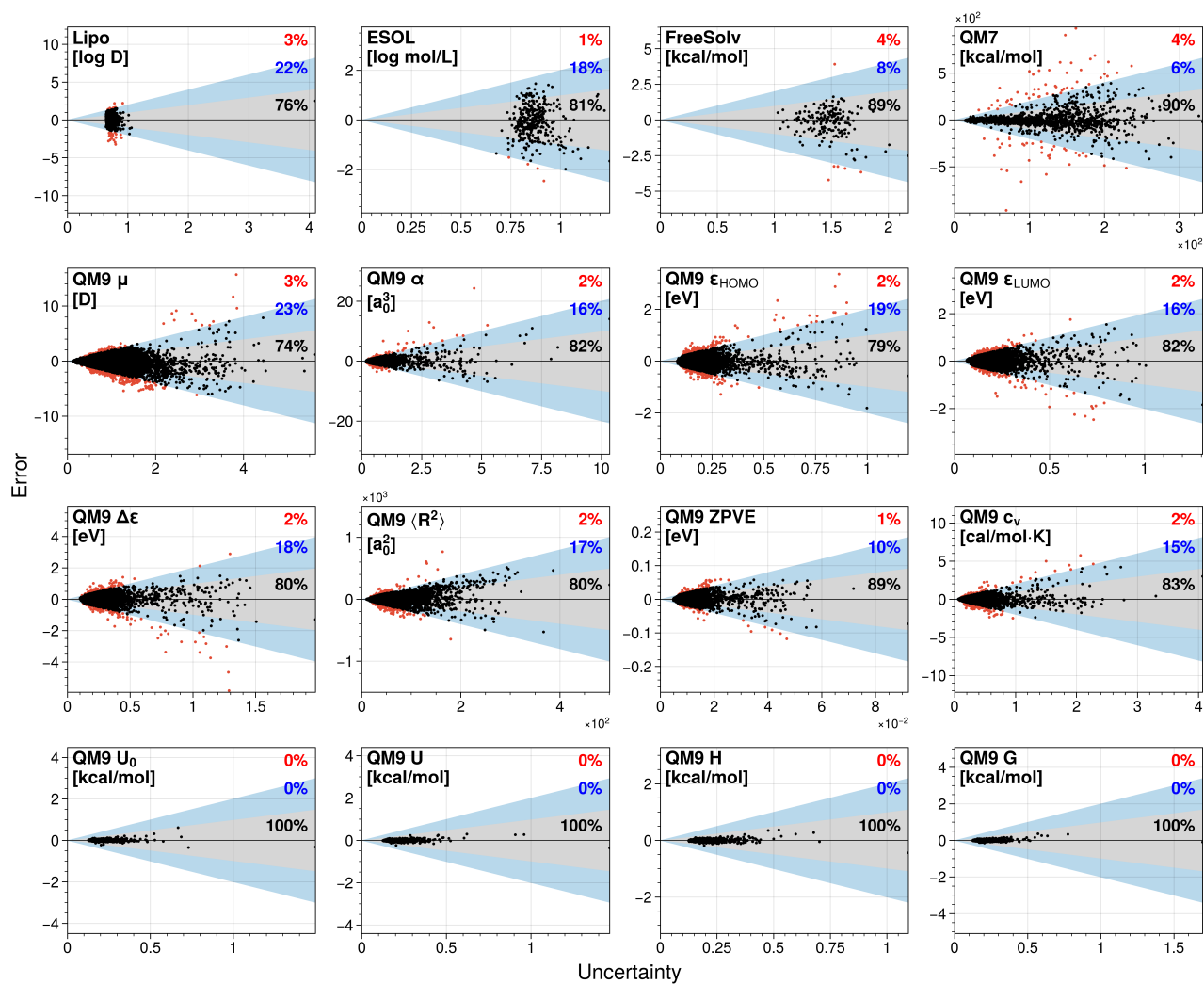


Figure S16: Prediction error vs uncertainty with a random seed of 1.

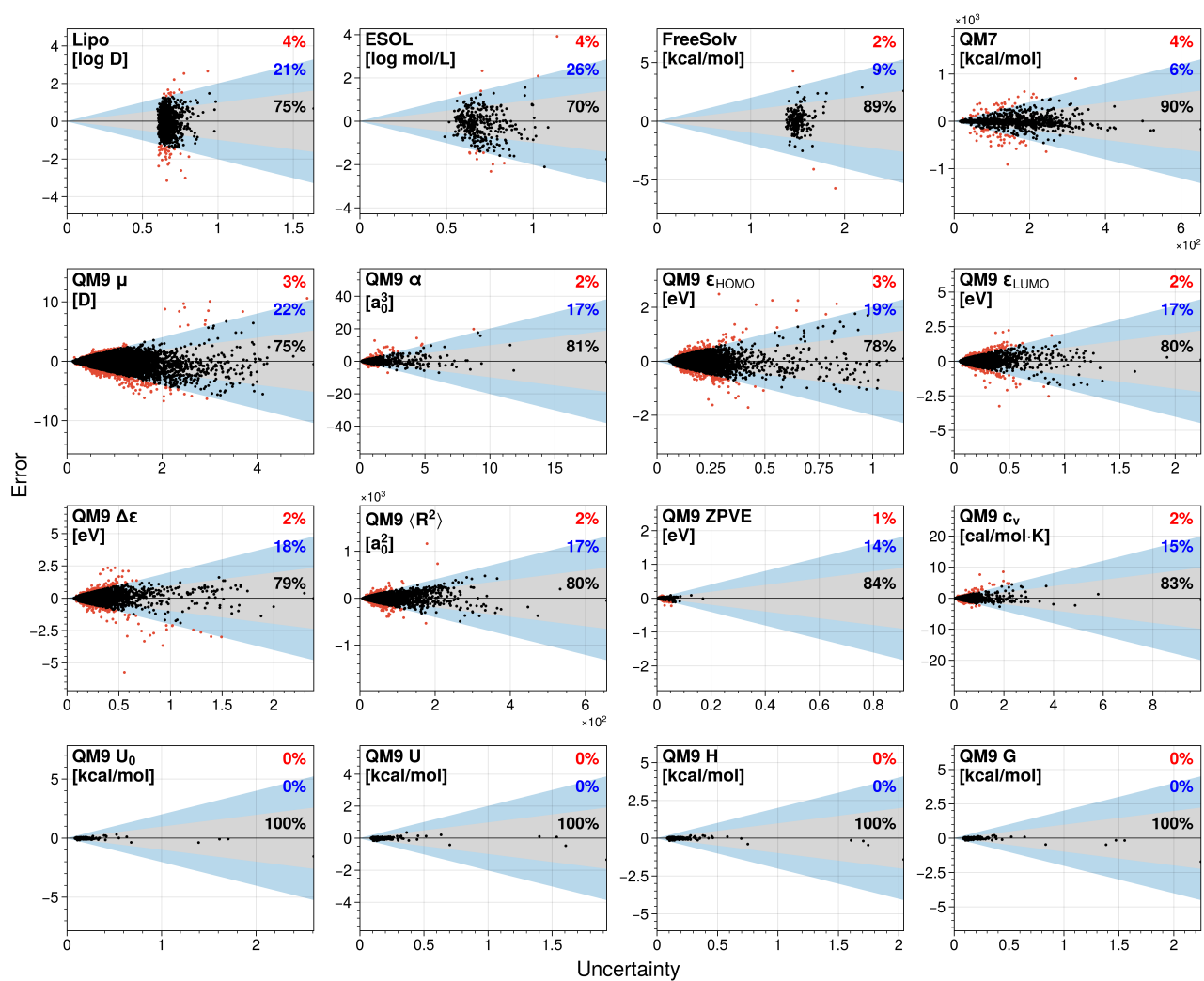


Figure S17: Prediction error vs uncertainty with a random seed of 2.

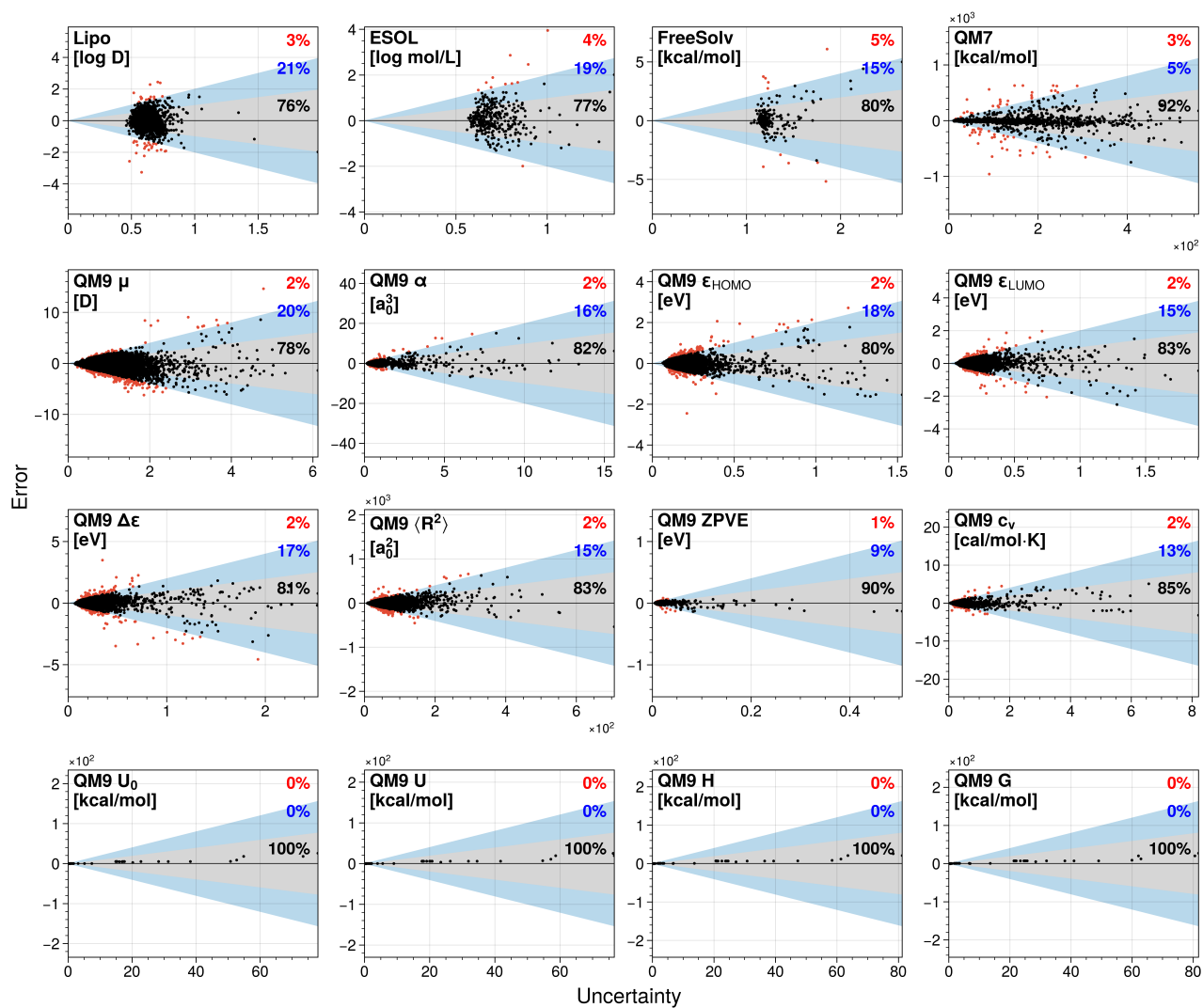


Figure S18: Prediction error vs uncertainty with a random seed of 3.

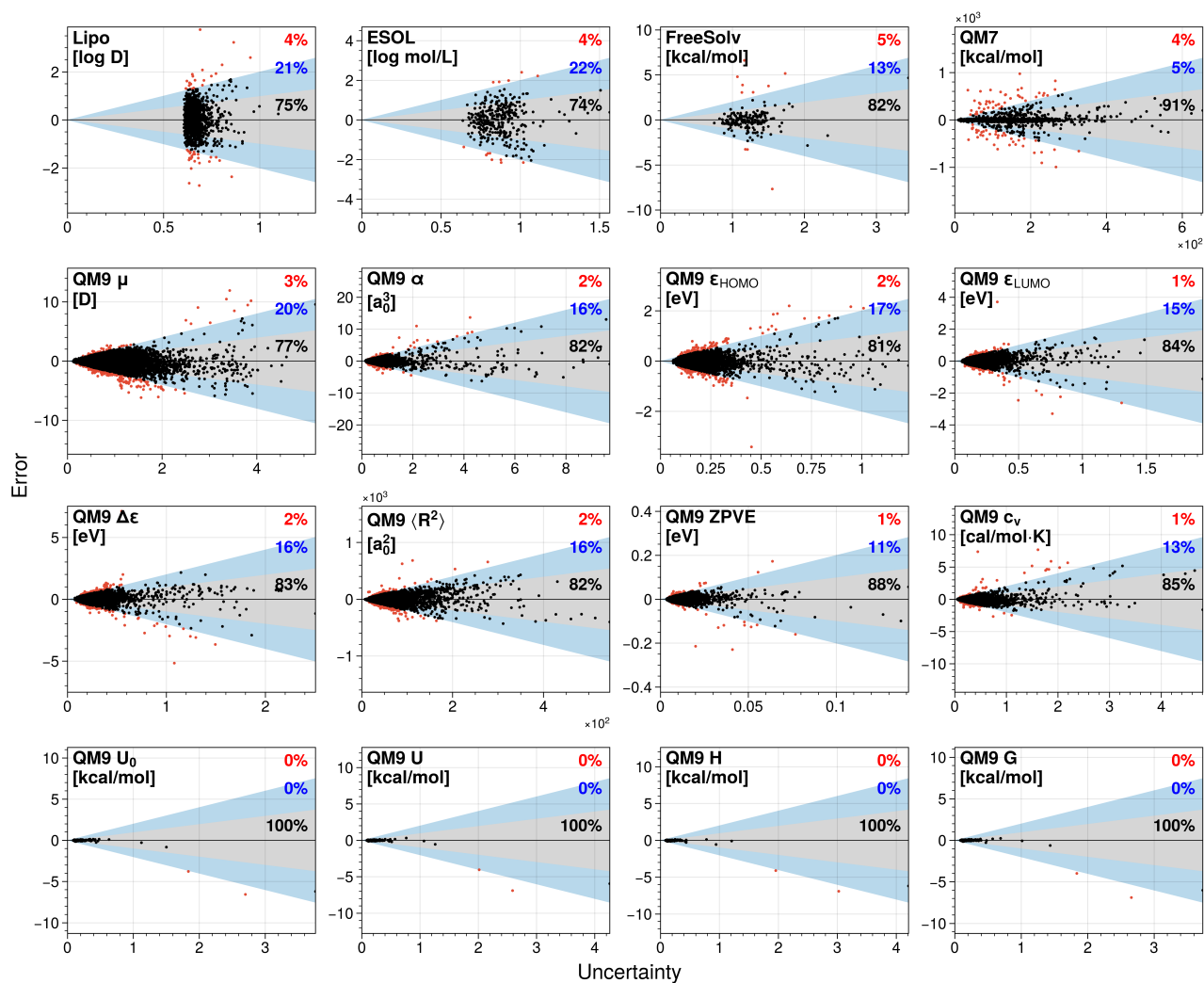


Figure S19: Prediction error vs uncertainty with a random seed of 4.

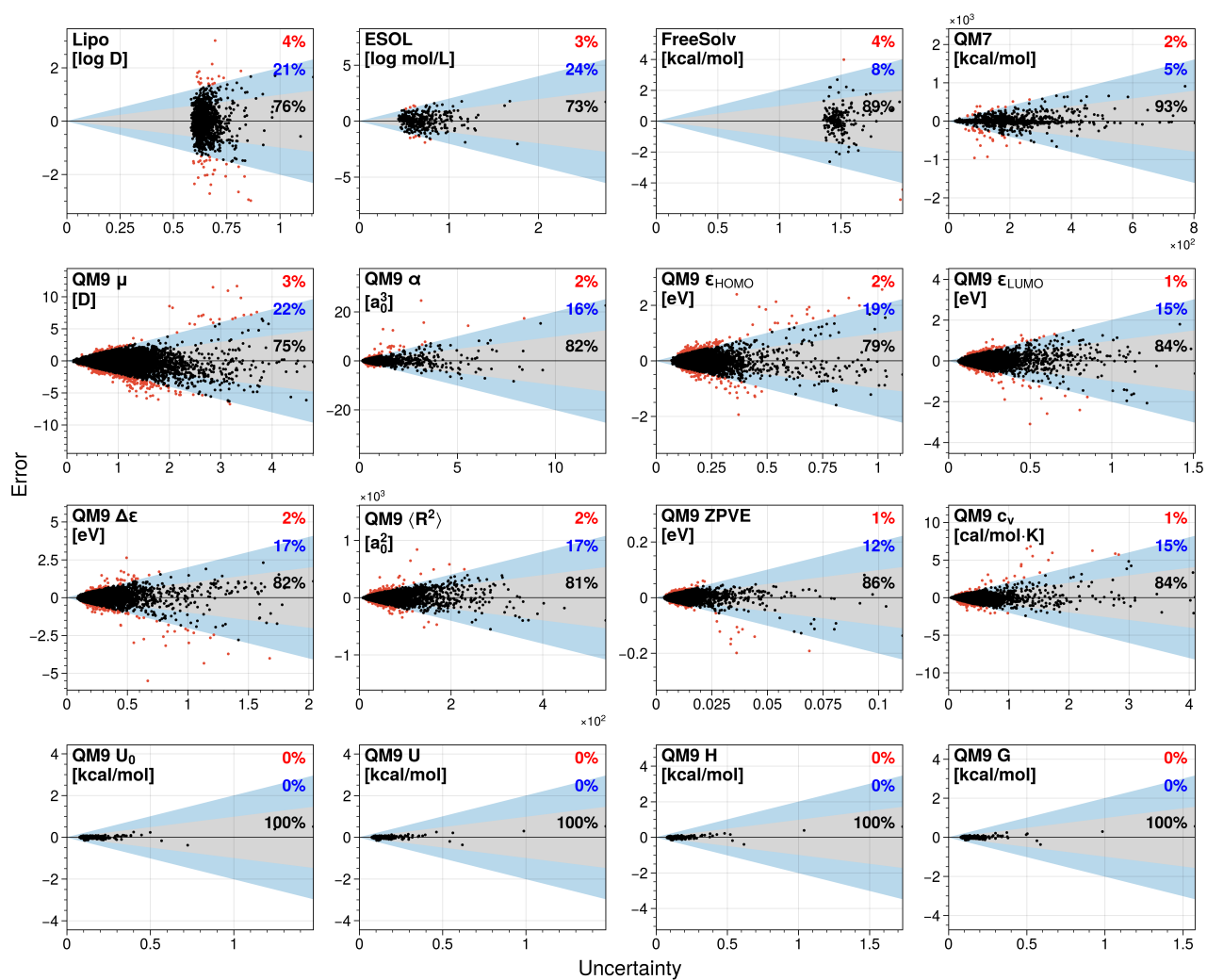


Figure S20: Prediction error vs uncertainty with a random seed of 5.

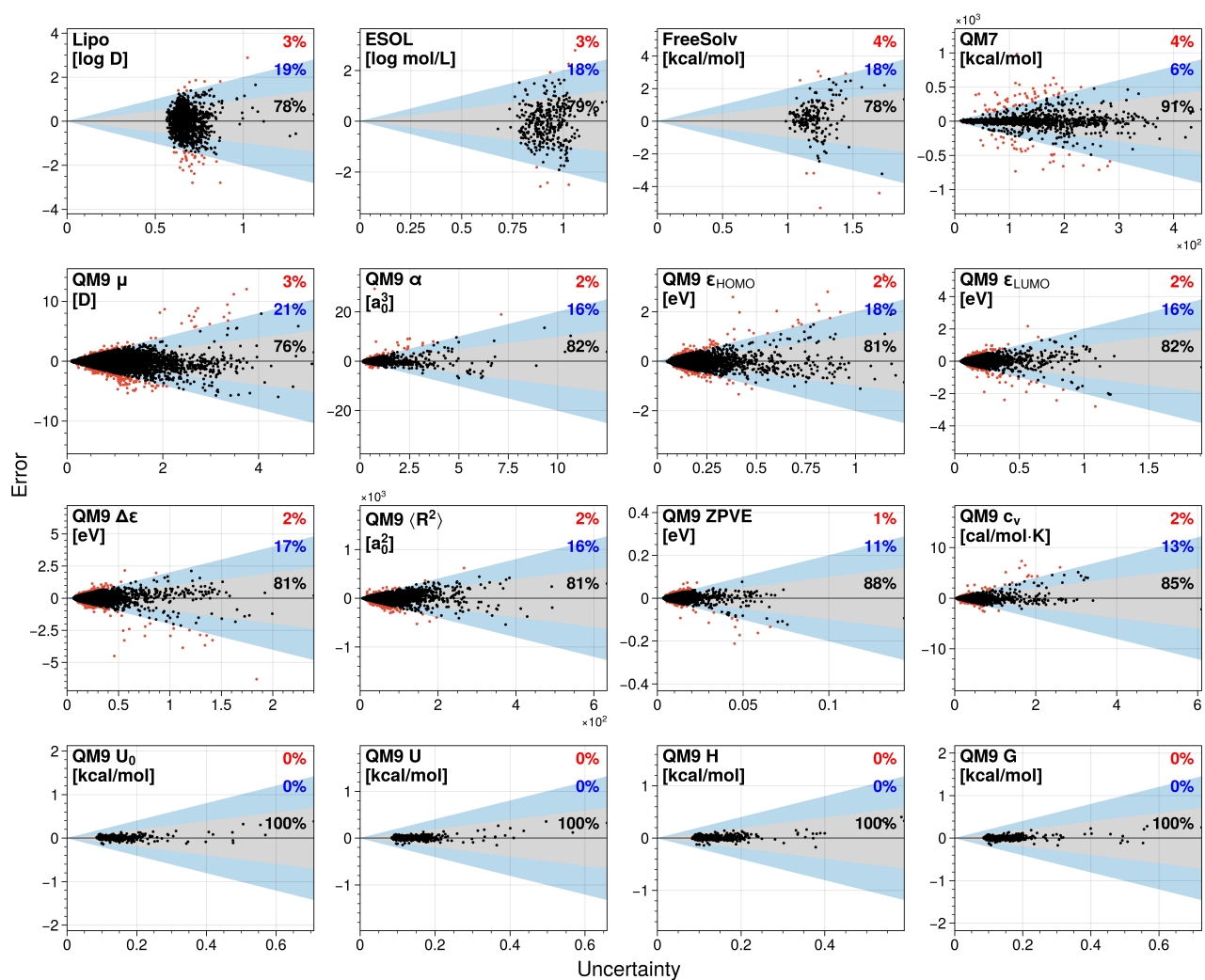


Figure S21: Prediction error vs uncertainty with a random seed of 6.

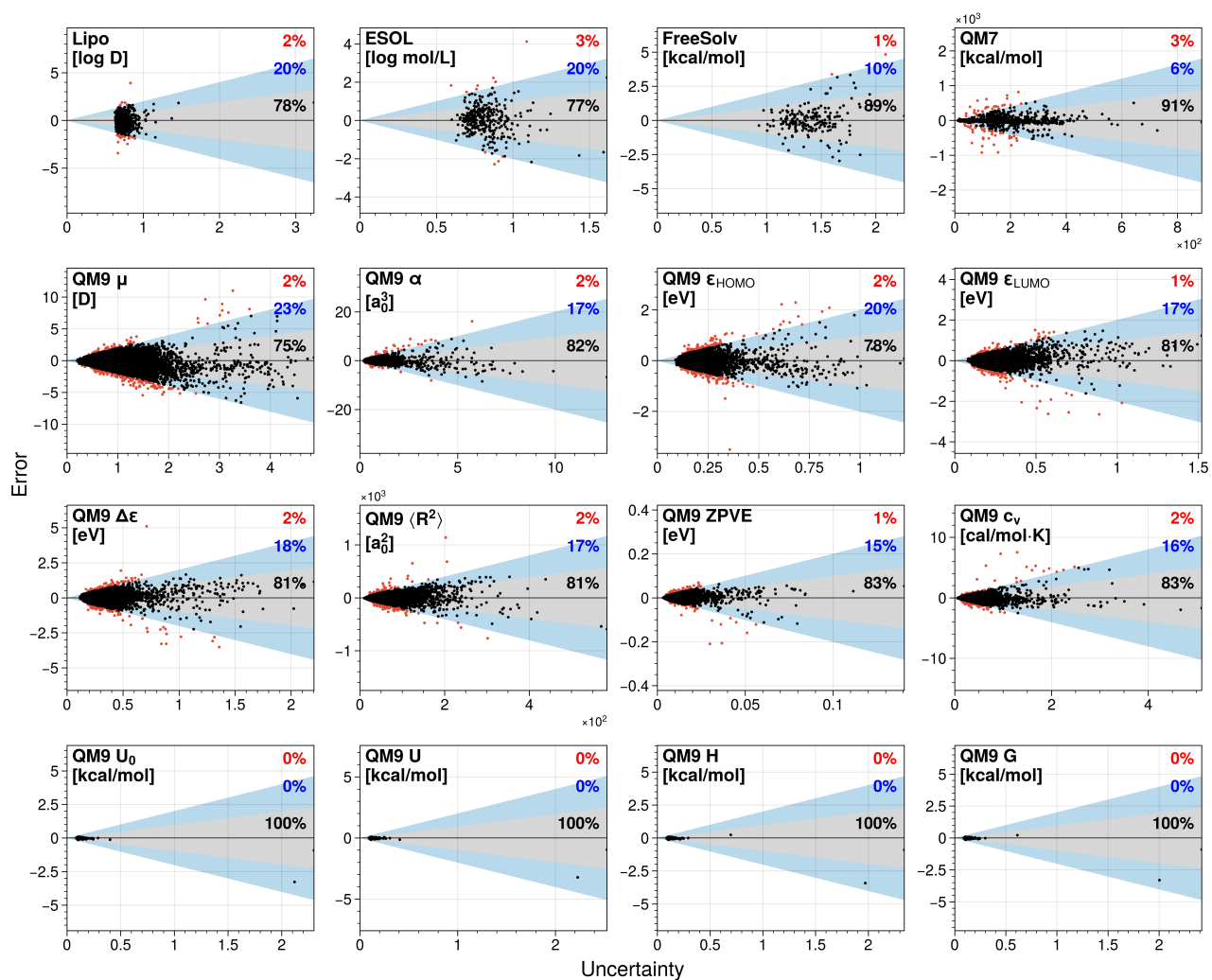


Figure S22: Prediction error vs uncertainty with a random seed of 7.

## S4 TSNE of AutoGNNUQ

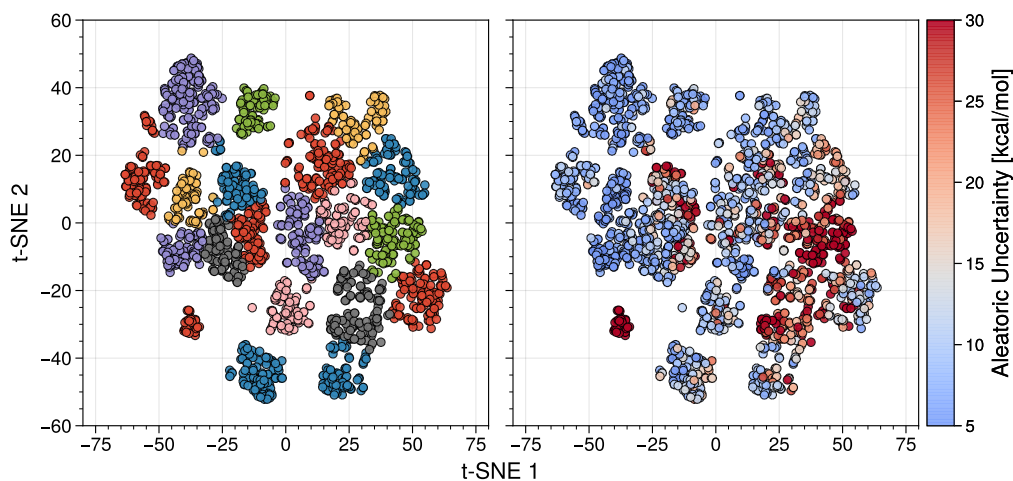


Figure S23: t-SNE plot of QM7 MACCS keys for random seed 0. (a) Clustering by k-means algorithm. (b) Aleatoric uncertainties associated with each molecule. Low uncertainties are represented in blue and high uncertainties in red.

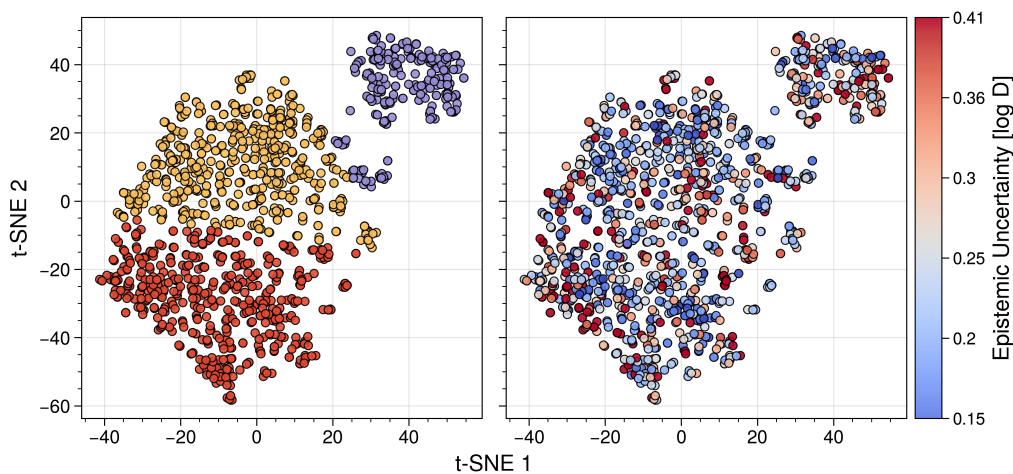


Figure S24: t-SNE plot of Lipo MACCS keys for random seed 0. (a) Clustering by k-means algorithm. (b) Epistemic uncertainties associated with each molecule. Low uncertainties are represented in blue and high uncertainties in red.



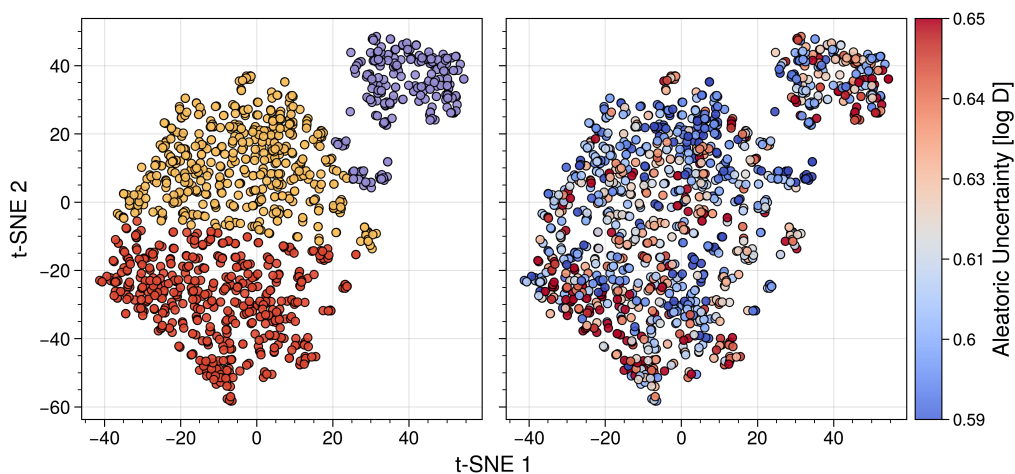


Figure S25: t-SNE plot of Lipo MACCS keys for random seed 0. (a) Clustering by k-means algorithm. (b) Aleatoric uncertainties associated with each molecule. Low uncertainties are represented in blue and high uncertainties in red.

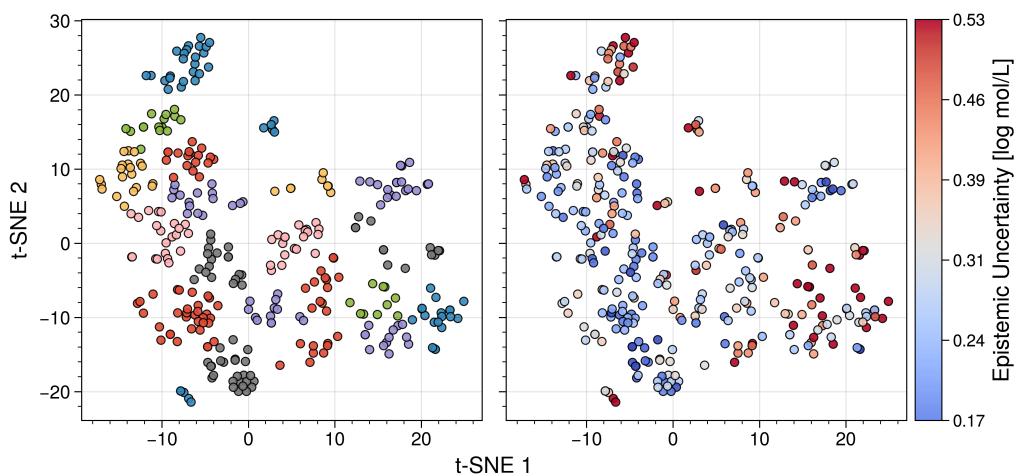


Figure S26: t-SNE plot of ESOL MACCS keys for random seed 0. (a) Clustering by k-means algorithm. (b) Epistemic uncertainties associated with each molecule. Low uncertainties are represented in blue and high uncertainties in red.

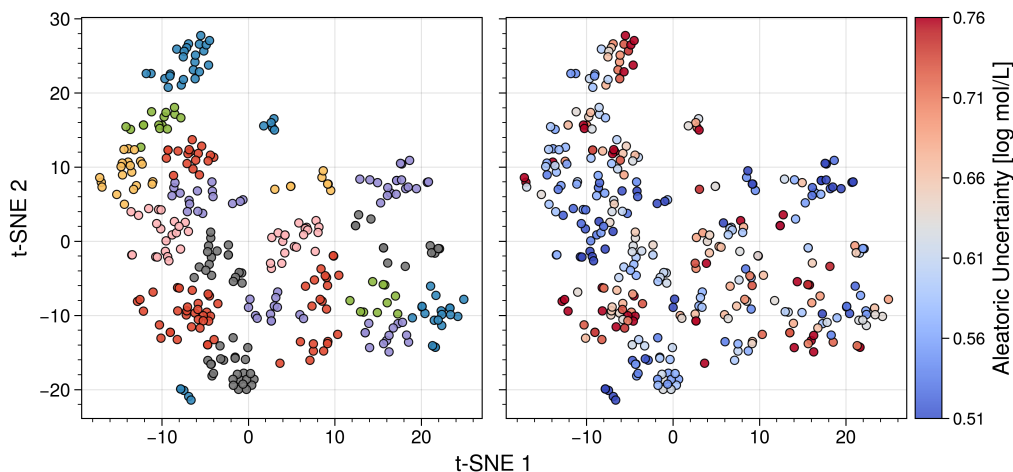


Figure S27: t-SNE plot of ESOL MACCS keys for random seed 0. (a) Clustering by k-means algorithm. (b) Aleatoric uncertainties associated with each molecule. Low uncertainties are represented in blue and high uncertainties in red.

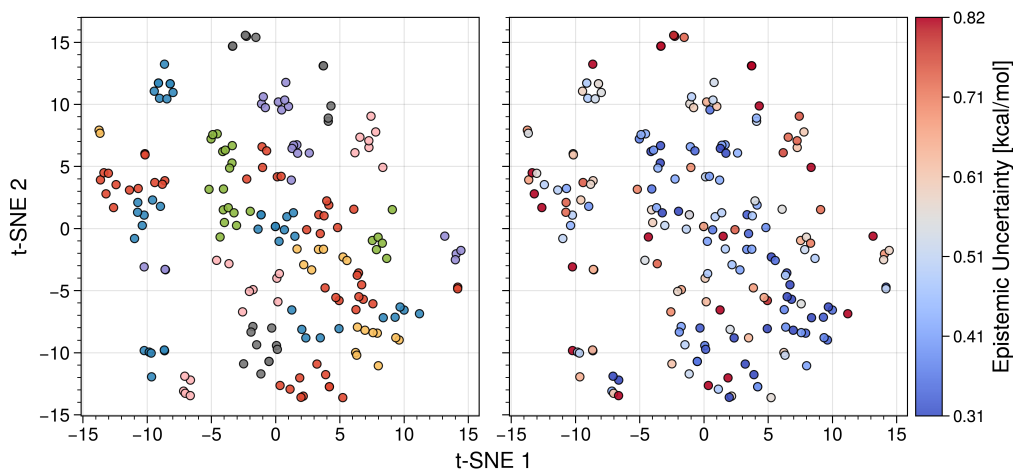


Figure S28: t-SNE plot of FreeSolv MACCS keys for random seed 0. (a) Clustering by k-means algorithm. (b) Epistemic uncertainties associated with each molecule. Low uncertainties are represented in blue and high uncertainties in red.

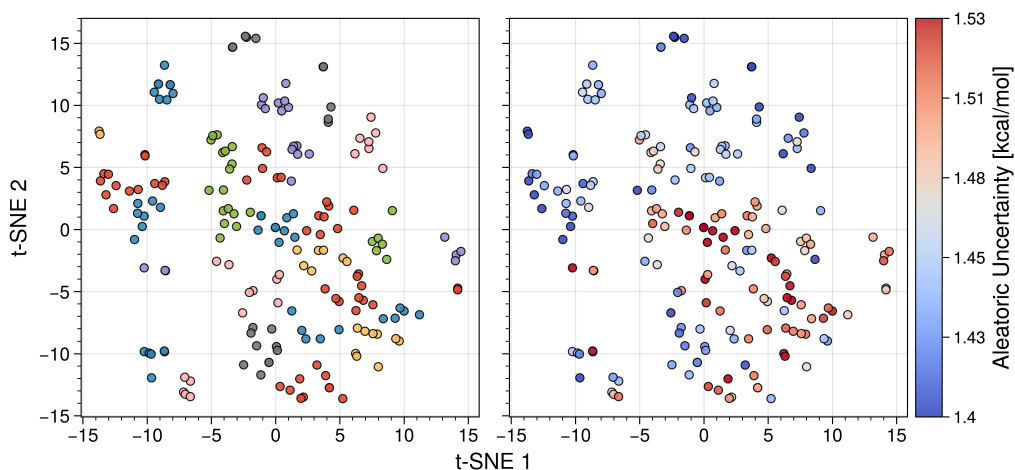


Figure S29: t-SNE plot of FreeSolv MACCS keys for random seed 0. (a) Clustering by k-means algorithm. (b) Aleatoric uncertainties associated with each molecule. Low uncertainties are represented in blue and high uncertainties in red.

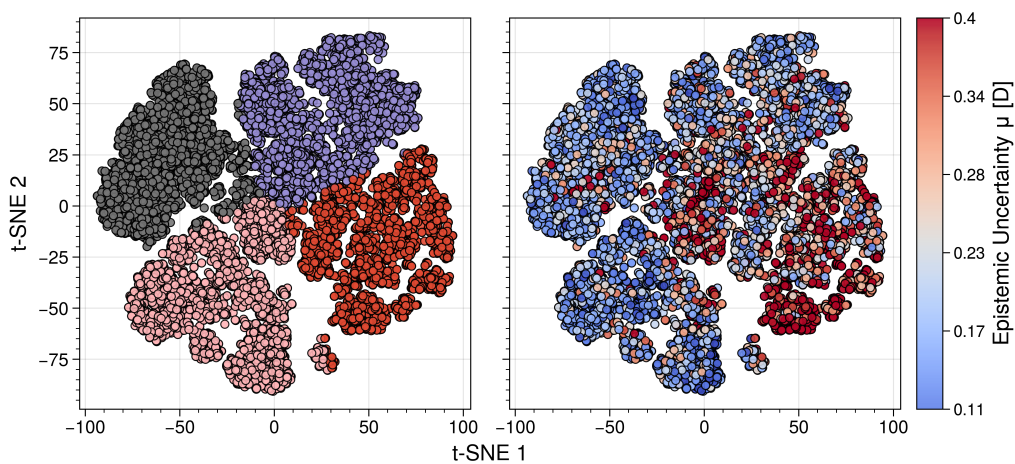


Figure S30: t-SNE plot of QM9  $\mu$  MACCS keys for random seed 0. (a) Clustering by k-means algorithm. (b) Epistemic uncertainties associated with each molecule. Low uncertainties are represented in blue and high uncertainties in red.

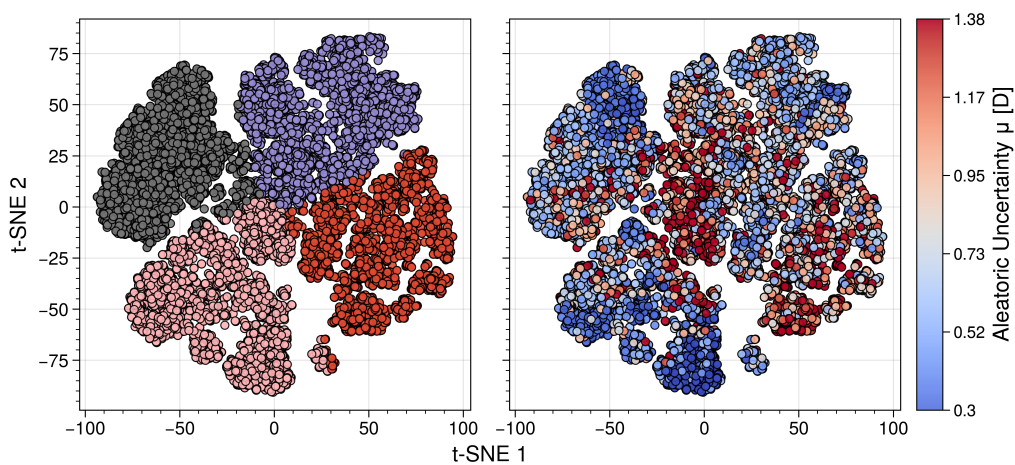


Figure S31: t-SNE plot of QM9  $\mu$  MACCS keys for random seed 0. (a) Clustering by k-means algorithm. (b) Aleatoric uncertainties associated with each molecule. Low uncertainties are represented in blue and high uncertainties in red.

## S5 Confidence-Based Recalibration

Table S2: Confidence-based calibration ratios. Mean value reported with standard deviation in parentheses.

Dataset	Property	AutoGNNUQ	AutoGNNUQ-Simple	MC Dropout	Random Ensemble
Lipo	logP	0.846 (0.035)	1.003 (0.087)	0.973 (0.080)	0.827 (0.028)
ESOL	$S_{\text{water}}$	0.844 (0.044)	0.894 (0.030)	1.003 (0.059)	0.652 (0.063)
FreeSolv	$\Delta G_{\text{hyd}}$	0.562 (0.067)	0.738 (0.131)	0.770 (0.150)	0.522 (0.050)
QM7	$\Delta H_{\text{atom}}^{\circ}$	0.276 (0.049)	0.252 (0.025)	0.380 (0.080)	0.358 (0.052)
QM9	$\mu$	0.847 (0.022)	0.919 (0.070)	0.874 (0.030)	0.799 (0.037)
	$\alpha$	0.721 (0.015)	0.730 (0.038)	0.507 (0.092)	0.441 (0.044)
	$\epsilon_{\text{HOMO}}$	0.773 (0.020)	0.792 (0.050)	0.768 (0.021)	0.634 (0.023)
	$\epsilon_{\text{LUMO}}$	0.709 (0.020)	0.726 (0.045)	0.669 (0.037)	0.546 (0.041)
	$\Delta\epsilon$	0.734 (0.020)	0.738 (0.050)	0.706 (0.024)	0.581 (0.039)
	$\langle R^2 \rangle$	0.732 (0.019)	0.751 (0.044)	0.646 (0.050)	0.505 (0.038)
	ZPVE	0.617 (0.048)	0.653 (0.069)	0.514 (0.115)	0.429 (0.082)
	$c_v$	0.670 (0.023)	0.695 (0.055)	0.566 (0.033)	0.465 (0.061)
	$U_0$	0.126 (0.035)	0.117 (0.022)	0.305 (0.046)	0.257 (0.063)
	$U$	0.112 (0.013)	0.116 (0.024)	0.305 (0.046)	0.255 (0.067)
	$H$	0.134 (0.022)	0.107 (0.013)	0.305 (0.046)	0.255 (0.067)
	$G$	0.114 (0.020)	0.124 (0.023)	0.305 (0.046)	0.256 (0.064)

Table S3: Confidence-based calibration metrics of random ensemble before and after recalibration based on the validation dataset. Mean value reported with standard deviation in parentheses.

Dataset	Property	ECE		MCE		MCA	
		Before	After	Before	After	Before	After
Lipo	logP	0.063 (0.008)	<b>0.009 (0.005)</b>	0.106 (0.013)	<b>0.023 (0.008)</b>	0.063 (0.008)	<b>0.009 (0.005)</b>
ESOL	$S_{\text{water}}$	0.135 (0.018)	<b>0.021 (0.013)</b>	0.223 (0.032)	<b>0.046 (0.019)</b>	0.136 (0.018)	<b>0.022 (0.013)</b>
FreeSolv	$\Delta G_{\text{hyd}}$	0.185 (0.015)	<b>0.035 (0.018)</b>	0.317 (0.037)	<b>0.079 (0.029)</b>	0.187 (0.015)	<b>0.035 (0.018)</b>
QM7	$\Delta H_{\text{atom}}^{\circ}$	0.241 (0.022)	<b>0.020 (0.003)</b>	0.442 (0.041)	<b>0.081 (0.014)</b>	0.243 (0.022)	<b>0.020 (0.003)</b>
QM9	$\mu$	0.072 (0.012)	<b>0.006 (0.001)</b>	0.114 (0.022)	<b>0.016 (0.003)</b>	0.073 (0.012)	<b>0.006 (0.001)</b>
	$\alpha$	0.242 (0.021)	<b>0.022 (0.003)</b>	0.435 (0.040)	<b>0.063 (0.011)</b>	0.245 (0.021)	<b>0.023 (0.003)</b>
	$\epsilon_{\text{HOMO}}$	0.141 (0.011)	<b>0.005 (0.001)</b>	0.224 (0.018)	<b>0.014 (0.003)</b>	0.142 (0.011)	<b>0.005 (0.001)</b>
	$\epsilon_{\text{LUMO}}$	0.194 (0.017)	<b>0.019 (0.005)</b>	0.345 (0.032)	<b>0.063 (0.021)</b>	0.196 (0.018)	<b>0.019 (0.005)</b>
	$\Delta\epsilon$	0.168 (0.020)	<b>0.017 (0.006)</b>	0.290 (0.032)	<b>0.038 (0.012)</b>	0.169 (0.020)	<b>0.017 (0.006)</b>
	$\langle R^2 \rangle$	0.203 (0.017)	<b>0.018 (0.007)</b>	0.356 (0.024)	<b>0.039 (0.015)</b>	0.205 (0.017)	<b>0.019 (0.007)</b>
	ZPVE	0.263 (0.043)	<b>0.044 (0.008)</b>	0.525 (0.086)	<b>0.138 (0.020)</b>	0.266 (0.043)	<b>0.044 (0.008)</b>
	$c_v$	0.237 (0.027)	<b>0.027 (0.007)</b>	0.438 (0.046)	<b>0.084 (0.025)</b>	0.240 (0.028)	<b>0.028 (0.007)</b>
	$U_0$	0.344 (0.031)	<b>0.042 (0.014)</b>	0.656 (0.043)	<b>0.106 (0.023)</b>	0.348 (0.032)	<b>0.043 (0.014)</b>
	$U$	0.345 (0.033)	<b>0.044 (0.010)</b>	0.656 (0.043)	<b>0.106 (0.024)</b>	0.349 (0.033)	<b>0.044 (0.010)</b>
	$H$	0.345 (0.033)	<b>0.044 (0.010)</b>	0.656 (0.044)	<b>0.107 (0.022)</b>	0.349 (0.033)	<b>0.044 (0.011)</b>
	$G$	0.345 (0.032)	<b>0.044 (0.011)</b>	0.657 (0.044)	<b>0.105 (0.025)</b>	0.348 (0.032)	<b>0.044 (0.011)</b>

Table S4: Confidence-based calibration metrics of MC dropout before and after recalibration based on the validation dataset. Mean value reported with standard deviation in parentheses.

Dataset	Property	ECE		MCE		MCA	
		Before	After	Before	After	Before	After
Lipo	logP	0.020 (0.010)	<b>0.014 (0.008)</b>	0.043 (0.016)	<b>0.031 (0.011)</b>	0.020 (0.010)	<b>0.014 (0.008)</b>
ESOL	$S_{\text{water}}$	0.017 (0.005)	<b>0.019 (0.005)</b>	0.043 (0.011)	<b>0.048 (0.009)</b>	0.017 (0.005)	<b>0.020 (0.005)</b>
FreeSolv	$\Delta G_{\text{hyd}}$	0.090 (0.049)	<b>0.035 (0.009)</b>	0.174 (0.080)	<b>0.094 (0.026)</b>	0.091 (0.049)	<b>0.035 (0.009)</b>
QM7	$\Delta H_{\text{atom}}^{\circ}$	0.215 (0.030)	<b>0.027 (0.007)</b>	0.383 (0.059)	<b>0.110 (0.022)</b>	0.217 (0.030)	<b>0.028 (0.007)</b>
QM9	$\mu$	0.042 (0.009)	<b>0.005 (0.002)</b>	0.067 (0.014)	<b>0.012 (0.006)</b>	0.043 (0.009)	<b>0.005 (0.002)</b>
	$\alpha$	0.200 (0.040)	<b>0.004 (0.001)</b>	0.322 (0.068)	<b>0.010 (0.002)</b>	0.202 (0.041)	<b>0.004 (0.001)</b>
	$\epsilon_{\text{HOMO}}$	0.081 (0.007)	<b>0.006 (0.002)</b>	0.128 (0.011)	<b>0.016 (0.004)</b>	0.082 (0.007)	<b>0.006 (0.002)</b>
	$\epsilon_{\text{LUMO}}$	0.120 (0.012)	<b>0.006 (0.001)</b>	0.189 (0.019)	<b>0.018 (0.005)</b>	0.121 (0.012)	<b>0.006 (0.001)</b>
	$\Delta\epsilon$	0.103 (0.011)	<b>0.006 (0.001)</b>	0.162 (0.016)	<b>0.018 (0.003)</b>	0.104 (0.011)	<b>0.006 (0.001)</b>
	$\langle R^2 \rangle$	0.130 (0.021)	<b>0.006 (0.002)</b>	0.207 (0.033)	<b>0.016 (0.004)</b>	0.131 (0.021)	<b>0.006 (0.002)</b>
	ZPVE	0.191 (0.057)	<b>0.008 (0.005)</b>	0.306 (0.095)	<b>0.023 (0.013)</b>	0.193 (0.058)	<b>0.008 (0.005)</b>
	$c_v$	0.166 (0.013)	<b>0.007 (0.001)</b>	0.262 (0.023)	<b>0.019 (0.004)</b>	0.168 (0.014)	<b>0.007 (0.002)</b>
	$U_0$	0.299 (0.042)	<b>0.011 (0.006)</b>	0.504 (0.059)	<b>0.037 (0.035)</b>	0.302 (0.042)	<b>0.011 (0.006)</b>
	$U$	0.299 (0.042)	<b>0.011 (0.006)</b>	0.504 (0.059)	<b>0.037 (0.036)</b>	0.302 (0.042)	<b>0.011 (0.006)</b>
	$H$	0.299 (0.042)	<b>0.011 (0.006)</b>	0.504 (0.059)	<b>0.037 (0.035)</b>	0.302 (0.042)	<b>0.011 (0.006)</b>
	$G$	0.299 (0.042)	<b>0.011 (0.006)</b>	0.504 (0.059)	<b>0.037 (0.036)</b>	0.302 (0.042)	<b>0.011 (0.006)</b>

Table S5: Confidence-based calibration metrics of `AutoGNNUQ-Simple` before and after recalibration based on the validation dataset. Mean value reported with standard deviation in parentheses.

Dataset	Property	ECE		MCE		MCA	
		Before	After	Before	After	Before	After
Lipo	logP	0.025 (0.011)	<b>0.012 (0.006)</b>	0.048 (0.017)	<b>0.026 (0.011)</b>	0.025 (0.011)	<b>0.012 (0.006)</b>
ESOL	$S_{\text{water}}$	0.040 (0.025)	<b>0.026 (0.013)</b>	0.082 (0.034)	<b>0.057 (0.019)</b>	0.041 (0.025)	<b>0.026 (0.013)</b>
FreeSolv	$\Delta G_{\text{hyd}}$	0.084 (0.051)	<b>0.037 (0.015)</b>	0.159 (0.071)	<b>0.079 (0.030)</b>	0.085 (0.051)	<b>0.037 (0.015)</b>
QM7	$\Delta H_{\text{atom}}^{\circ}$	0.271 (0.013)	<b>0.031 (0.006)</b>	0.488 (0.022)	<b>0.129 (0.011)</b>	0.274 (0.013)	<b>0.032 (0.006)</b>
QM9	$\mu$	0.032 (0.019)	<b>0.007 (0.003)</b>	0.058 (0.027)	<b>0.018 (0.010)</b>	0.032 (0.020)	<b>0.007 (0.004)</b>
	$\alpha$	0.094 (0.015)	<b>0.006 (0.002)</b>	0.147 (0.022)	<b>0.018 (0.005)</b>	0.095 (0.015)	<b>0.006 (0.002)</b>
	$\epsilon_{\text{HOMO}}$	0.071 (0.017)	<b>0.005 (0.002)</b>	0.114 (0.025)	<b>0.012 (0.003)</b>	0.072 (0.017)	<b>0.005 (0.002)</b>
	$\epsilon_{\text{LUMO}}$	0.097 (0.017)	<b>0.005 (0.002)</b>	0.153 (0.027)	<b>0.013 (0.005)</b>	0.098 (0.018)	<b>0.005 (0.002)</b>
	$\Delta\epsilon$	0.091 (0.016)	<b>0.007 (0.002)</b>	0.144 (0.026)	<b>0.017 (0.006)</b>	0.092 (0.017)	<b>0.007 (0.002)</b>
	$\langle R^2 \rangle$	0.085 (0.018)	<b>0.007 (0.003)</b>	0.130 (0.027)	<b>0.016 (0.005)</b>	0.086 (0.018)	<b>0.007 (0.003)</b>
	ZPVE	0.125 (0.029)	<b>0.007 (0.002)</b>	0.198 (0.046)	<b>0.023 (0.006)</b>	0.126 (0.029)	<b>0.007 (0.002)</b>
	$c_v$	0.104 (0.023)	<b>0.009 (0.002)</b>	0.164 (0.036)	<b>0.027 (0.007)</b>	0.105 (0.023)	<b>0.009 (0.003)</b>
	$U_0$	0.419 (0.013)	<b>0.006 (0.002)</b>	0.761 (0.029)	<b>0.017 (0.004)</b>	0.424 (0.013)	<b>0.006 (0.002)</b>
	$U$	0.420 (0.014)	<b>0.008 (0.005)</b>	0.762 (0.028)	<b>0.022 (0.008)</b>	0.424 (0.014)	<b>0.009 (0.006)</b>
	$H$	0.424 (0.008)	<b>0.008 (0.002)</b>	0.768 (0.023)	<b>0.024 (0.007)</b>	0.428 (0.008)	<b>0.008 (0.002)</b>
	$G$	0.416 (0.013)	<b>0.010 (0.008)</b>	0.757 (0.023)	<b>0.025 (0.013)</b>	0.420 (0.013)	<b>0.011 (0.008)</b>

## S6 Confidence Curve and Comparison with Evidential Uncertainty Quantification

Table S6: Model error (RMSE) at various confidence percentile cutoffs. Mean values are reported with the standard error of the mean in parentheses. The better result is bold and underlined. The random split ratio of `AutoGNNUQ` is 8:1:1 to be consistent with the benchmark evidential UQ [1].

Dataset	AutoGNNUQ	Benchmark	Dataset	AutoGNNUQ	Benchmark
Lipo	0.61 (0.01)	<b><u>0.55 (0.02)</u></b>	ESOL	0.73 (0.04)	<b><u>0.66 (0.02)</u></b>
	<b><u>0.49 (0.01)</u></b>	0.50 (0.01)		0.55 (0.03)	<b><u>0.44 (0.01)</u></b>
	<b><u>0.45 (0.02)</u></b>	0.51 (0.02)		0.47 (0.04)	<b><u>0.35 (0.02)</u></b>
	<b><u>0.46 (0.04)</u></b>	0.53 (0.03)		0.51 (0.06)	<b><u>0.28 (0.02)</u></b>
	<b><u>0.47 (0.06)</u></b>	0.50 (0.04)		0.53 (0.07)	<b><u>0.22 (0.02)</u></b>
	1.45 (0.12)	<b><u>0.96 (0.07)</u></b>		<b><u>107 (3)</u></b>	115 (3)
FreeSolv	0.86 (0.06)	<b><u>0.42 (0.04)</u></b>	QM7	<b><u>25 (4)</u></b>	39 (3)
	0.85 (0.08)	<b><u>0.26 (0.04)</u></b>		<b><u>7 (1)</u></b>	23 (4)
	0.85 (0.10)	<b><u>0.25 (0.08)</u></b>		<b><u>5 (1)</u></b>	10 (4)
	0.79 (0.17)	<b><u>0.28 (0.12)</u></b>		<b><u>5 (1)</u></b>	10 (4)



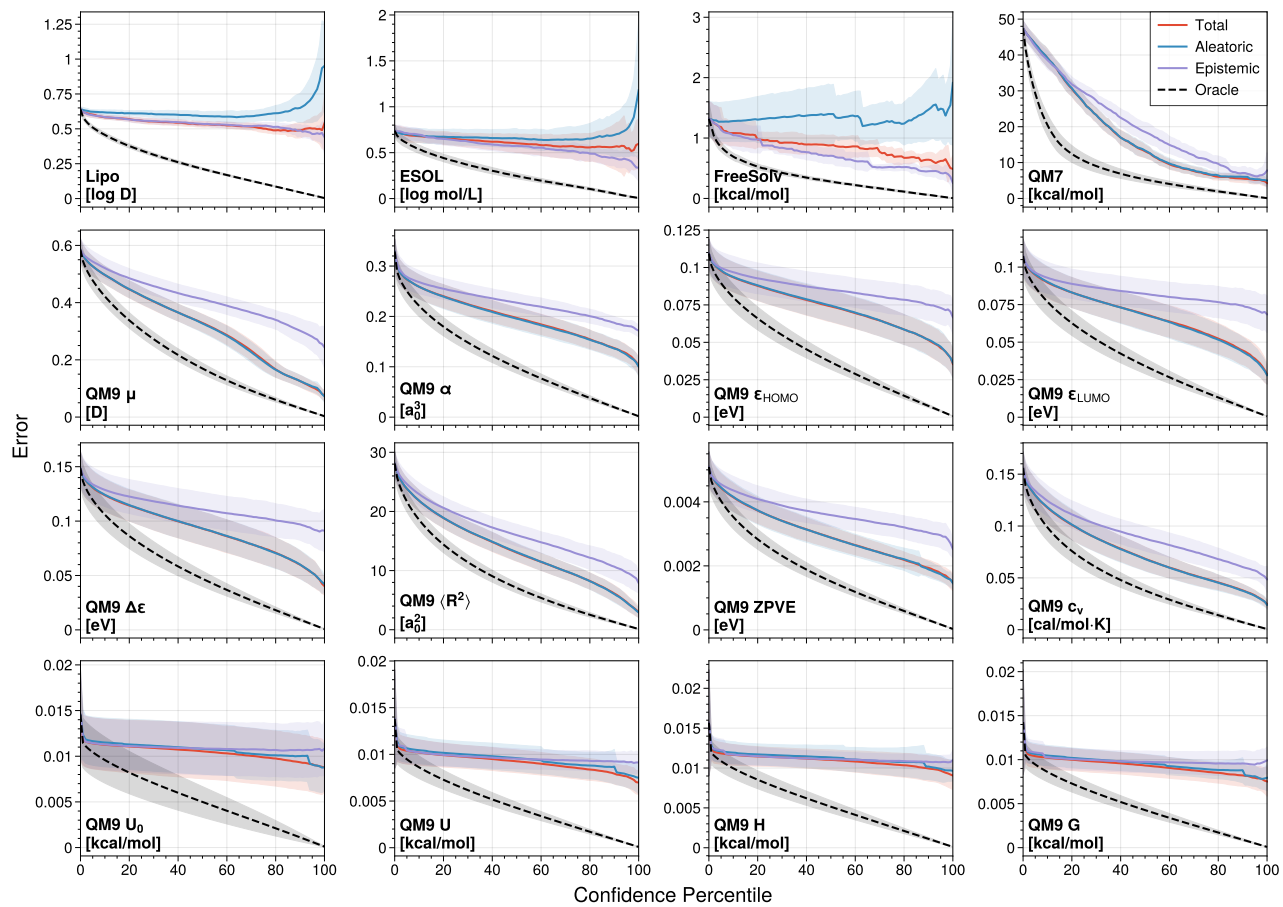


Figure S32: Confidence curve of AutoUQGNN of all test datasets across eight random seeds.

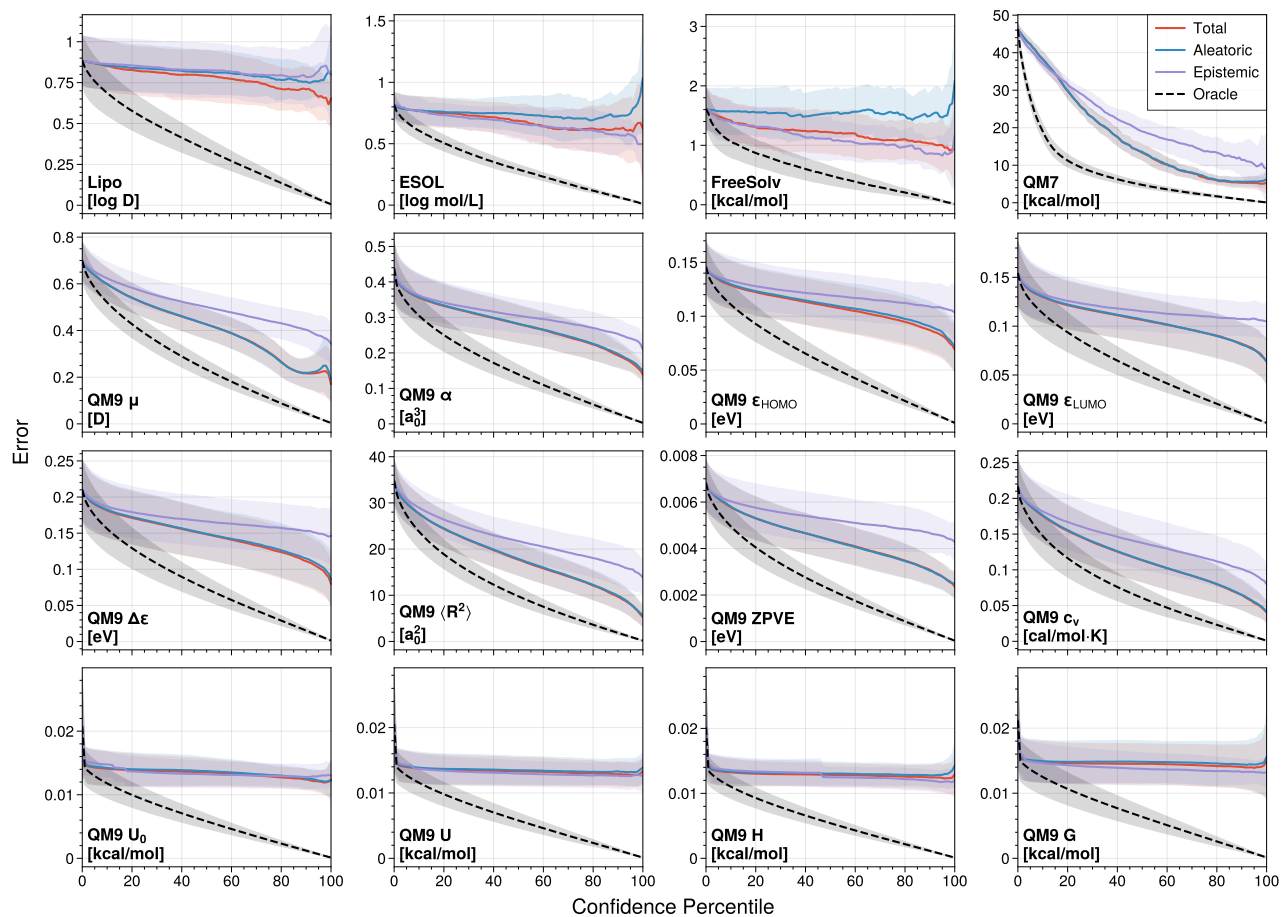


Figure S33: Confidence curve of AutoUQGNN-Simple of all test datasets across eight random seeds.

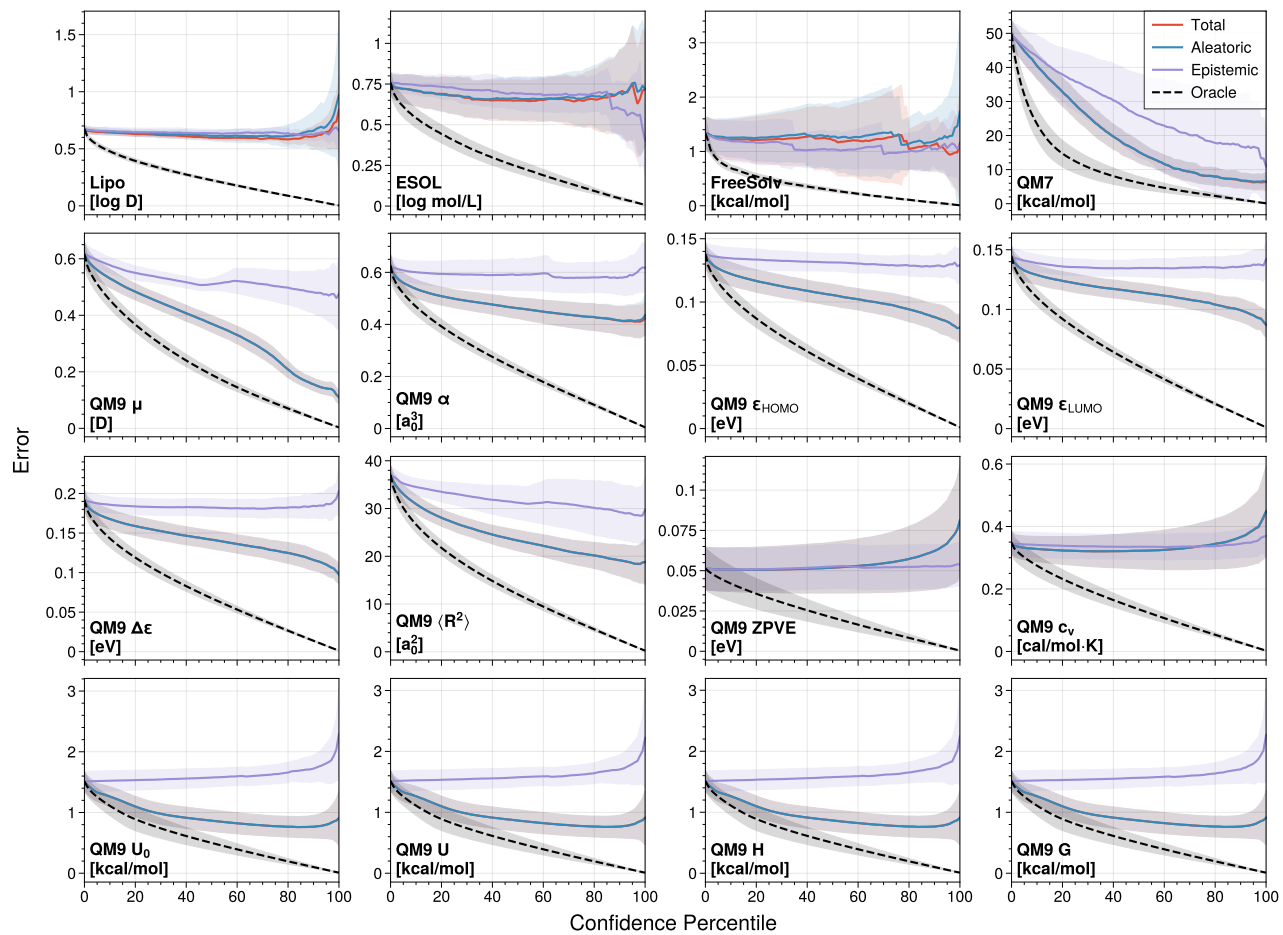


Figure S34: Confidence curve of MC dropout of all test datasets across eight random seeds.

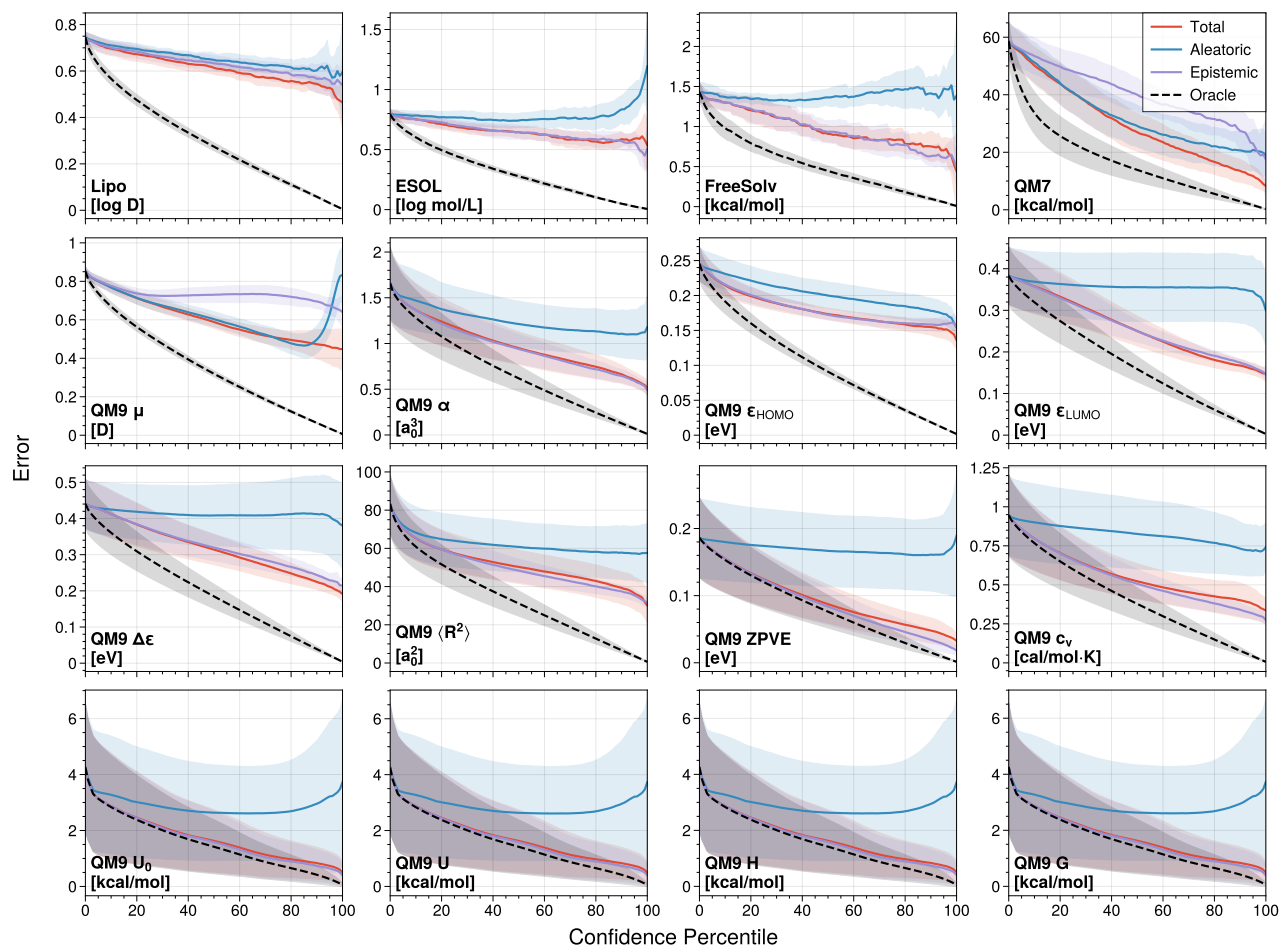


Figure S35: Confidence curve of random ensemble of all test datasets across eight random seeds.

Table S7: Confidence curve area under confidence-oracle error (AUCO). Mean value reported with standard deviation in parentheses. The best result is bold and underlined.

Dataset	Property	AutoGNNUQ	AutoGNNUQ-Simple	MC Dropout	Random Ensemble
Lipo	logP	<u>30.55 (1.75)</u>	42.06 (6.92)	37.64 (3.07)	32.16 (1.66)
ESOL	$S_{\text{water}}$	<u>34.00 (8.98)</u>	37.46 (10.96)	40.54 (11.45)	34.67 (2.71)
FreeSolv	$\Delta G_{\text{hyd}}$	53.13 (10.23)	66.41 (18.86)	86.05 (47.86)	<u>48.63 (9.66)</u>
QM7	$\Delta H_{\text{atom}}^{\circ}$	<u>939.45 (63.19)</u>	989.97 (86.88)	1016.44 (104.93)	1332.97 (191.13)
QM9	$\mu$	<u>11.77 (1.10)</u>	15.87 (2.70)	13.82 (1.11)	26.11 (2.35)
	$\alpha$	<u>8.82 (0.95)</u>	12.48 (2.22)	22.91 (2.56)	32.47 (4.82)
	$\epsilon_{\text{HOMO}}$	<u>3.29 (0.52)</u>	5.13 (1.14)	5.20 (0.48)	8.14 (0.47)
	$\epsilon_{\text{LUMO}}$	<u>2.97 (0.50)</u>	4.88 (1.09)	5.78 (0.47)	8.91 (0.86)
	$\Delta\epsilon$	<u>4.00 (0.72)</u>	6.81 (1.84)	6.89 (0.67)	12.38 (1.21)
	$\langle R^2 \rangle$	<u>496.94 (64.98)</u>	684.88 (145.02)	1092.04 (139.38)	1887.69 (293.54)
	ZPVE	<u>0.12 (0.01)</u>	0.19 (0.04)	3.27 (1.17)	1.45 (0.64)
	$c_v$	<u>2.76 (0.45)</u>	4.62 (1.33)	19.45 (4.55)	16.17 (2.95)
	$U_0$	<u>0.53 (0.14)</u>	0.72 (0.10)	39.02 (7.78)	20.86 (17.48)
	$U$	<u>0.47 (0.06)</u>	0.73 (0.10)	39.04 (7.81)	20.81 (17.27)
	$H$	<u>0.56 (0.07)</u>	0.72 (0.10)	39.02 (7.81)	20.80 (17.29)
	$G$	<u>0.48 (0.06)</u>	0.79 (0.16)	39.03 (7.80)	20.72 (17.10)

## S7 Loss Curves

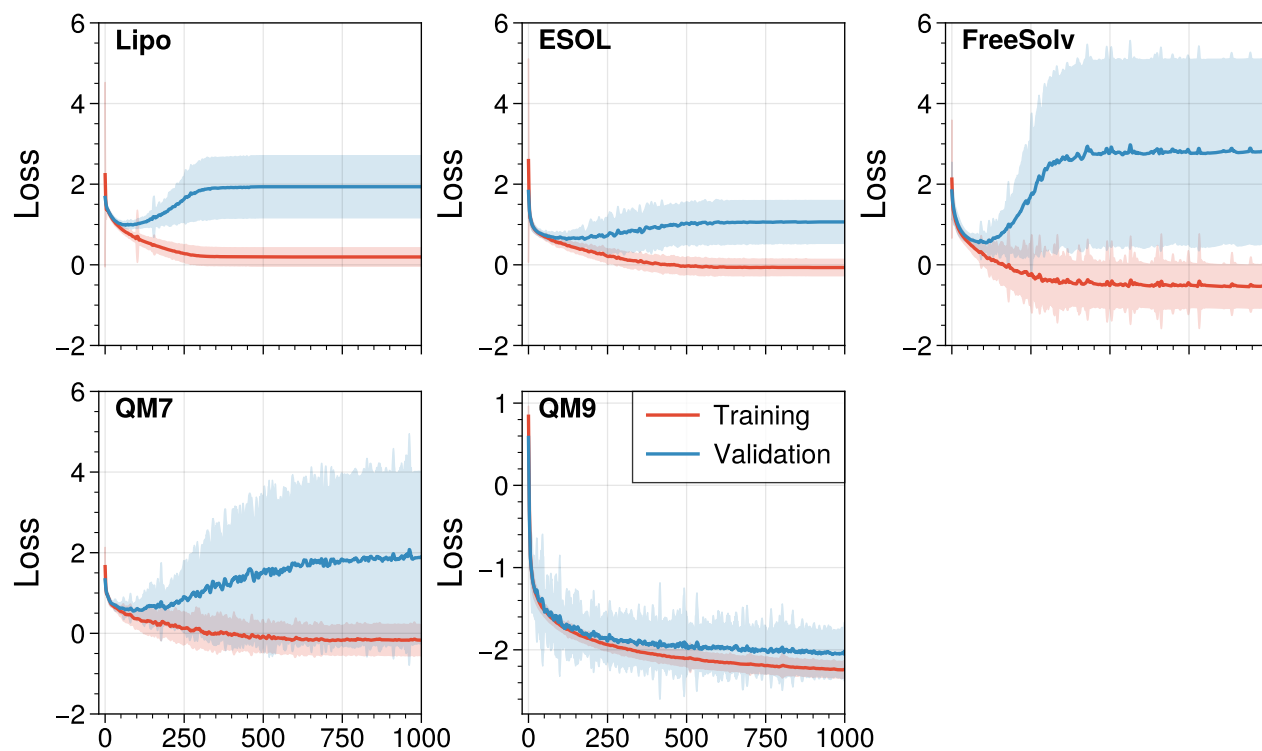


Figure S36: Loss curves of `AutoGNNUQ`. The solid line within the distribution represents the mean, while the shaded area signifies the standard deviation.

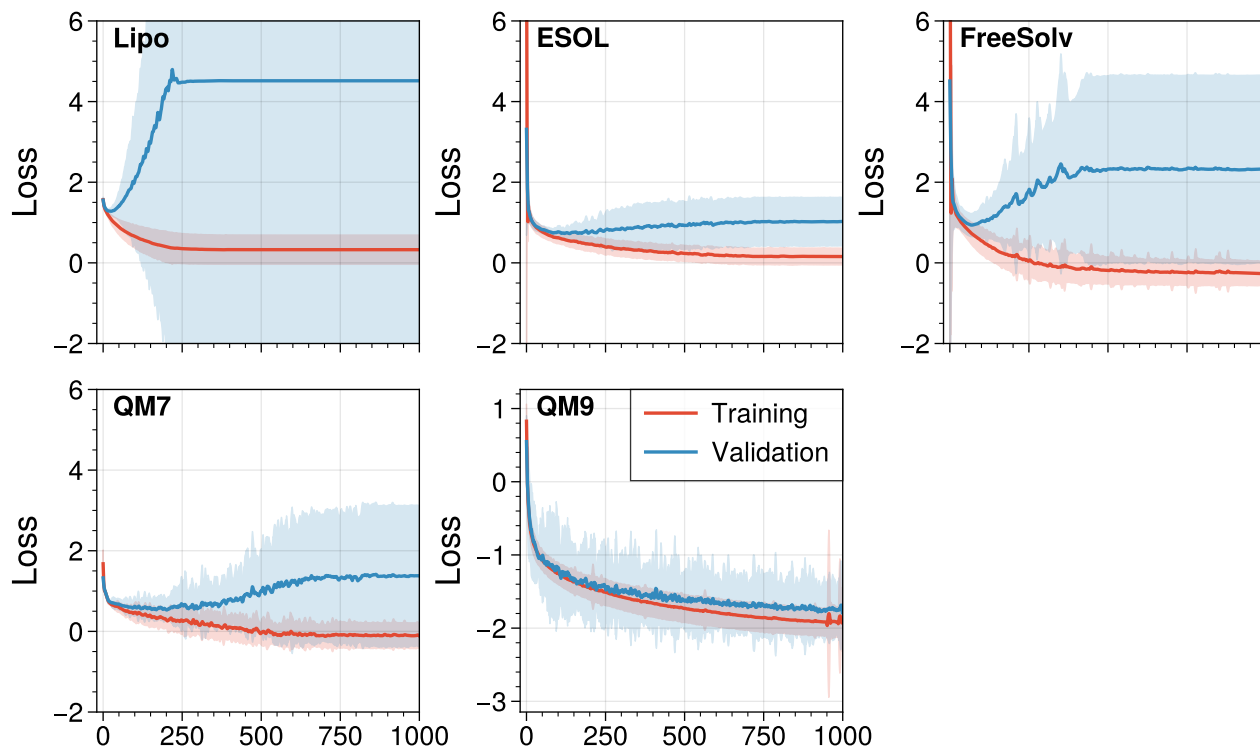


Figure S37: Loss curves of AutoGNNUQ-Simple. The solid line within the distribution represents the mean, while the shaded area signifies the standard deviation.

## S8 Out-of-Distribution Performance

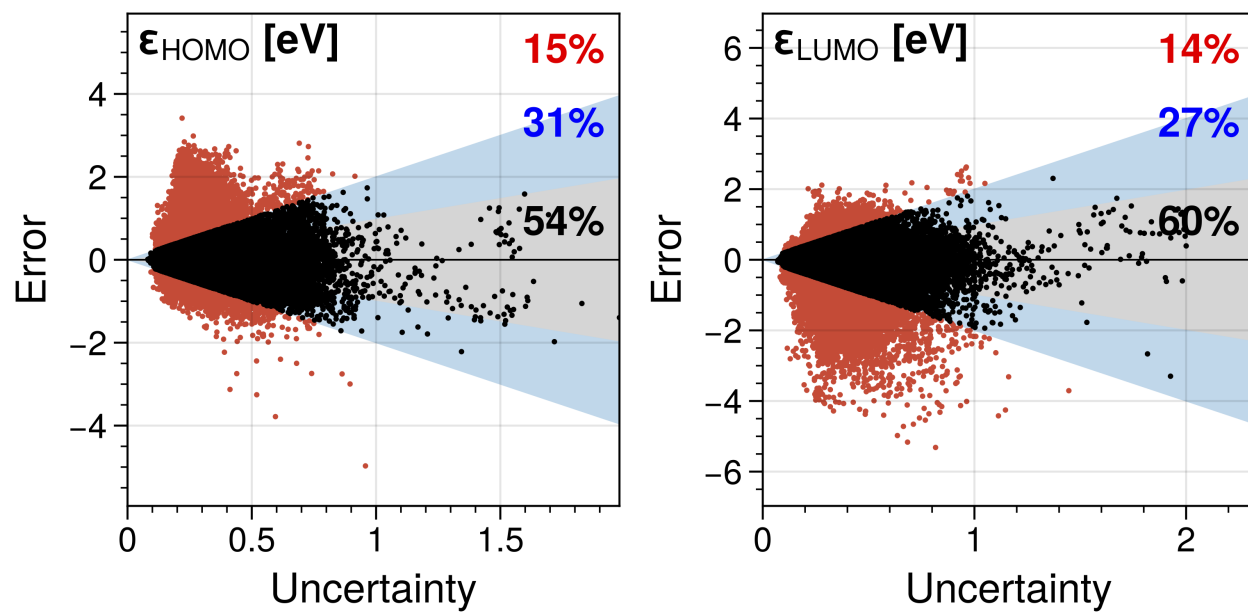


Figure S38: Prediction error vs uncertainty of PC9 using AutoGNNuQ trained on QM9.



## References

- [1] Soleimany, A. P. *et al.* Evidential deep learning for guided molecular property prediction and discovery. *ACS central science* **7**, 1356–1367 (2021).