

**Table 1: RMSE on three regression tasks**

	Morgan RF	Morgan XGB	Morgan DNN	MP graph NEF*	MP graph GCNN	MP graph Weave*	N-gram graph RF	N-gram graph XGB	N-gram graph DNN
delaney	1.311 $\pm 0.174$	1.110 $\pm 0.129$	1.231 $\pm 0.109$	<b>0.520</b> <b><math>\pm 0.137</math></b>	0.913 $\pm 0.061$	<b>0.460</b> <b><math>\pm 0.157</math></b>	0.773 $\pm 0.076$	0.700 $\pm 0.091$	<b>0.699</b> <b><math>\pm 0.046</math></b>
malaria	<b>1.028</b> <b><math>\pm 0.029</math></b>	<b>1.008</b> <b><math>\pm 0.031</math></b>	1.052 $\pm 0.051$	1.160 $\pm 0.059$	1.055 $\pm 0.042$	1.070 $\pm 0.118$	1.030 $\pm 0.023$	<b>1.010</b> <b><math>\pm 0.026</math></b>	1.119 $\pm 0.027$
cep	1.642 $\pm 0.026$	1.410 $\pm 0.040$	1.477 $\pm 0.042$	1.430 $\pm 0.176$	<b>1.184</b> <b><math>\pm 0.061</math></b>	<b>1.100</b> <b><math>\pm 0.118</math></b>	1.379 $\pm 0.013$	<b>1.290</b> <b><math>\pm 0.026</math></b>	1.365 $\pm 0.034$

- Mean and standard deviation on 5-fold cross validation.
- Top three results are **bolded**.
- Baseline results (marked in \*) are from [1][2].
- For model **Morgan DNN**, we were using the results from [1] in the submitted version. Then we got a hyperparameter set that performs better on all three tasks after submission, so we replace it here. (It was 1.4, 1.13, 2.0 on three tasks.)
- MP graph is short for Message-passing graph.

[1] Duvenaud, D. K.; Maclaurin, D.; Iparraguirre, J.; Bombarell, R.; Hirzel, T.; Aspuru-Guzik, A.; and Adams, R. P. 2015. Convolutional Networks on Graphs for Learning Molecular Fingerprints. 2224–2232.

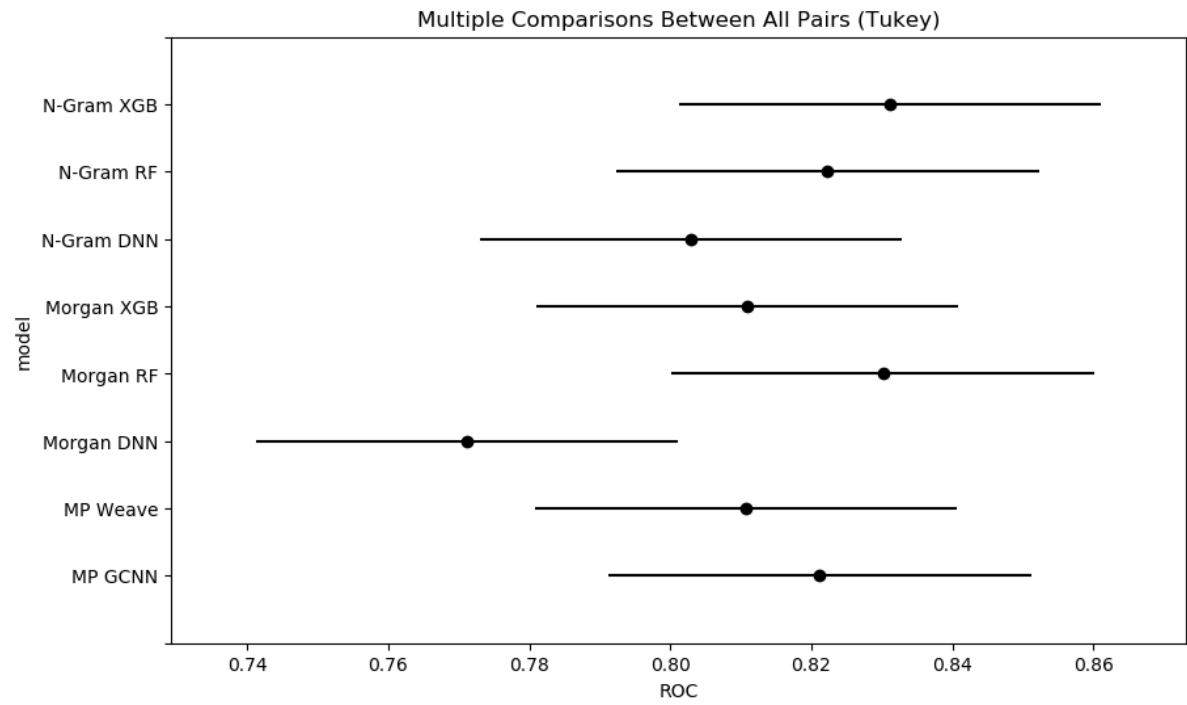
[2] Kusner, M. J.; Paige, B.; and Hernández-Lobato, J. M. 2017. Grammar variational autoencoder. arXiv preprint arXiv:1703.01925.

**Table 4: AUC[ROC] on Tox21**

	Morgan RF	Morgan XGB	Morgan DNN	MP graph GCNN	MP graph Weave	N-gram graph RF	N-gram graph XGB	N-gram graph DNN
NR-AR	0.787 $\pm 0.065$	0.777 $\pm 0.062$	0.756 $\pm 0.057$	0.793 $\pm 0.068$	<b>0.796</b> <b><math>\pm 0.057</math></b>	<b>0.802</b> <b><math>\pm 0.082</math></b>	0.790 $\pm 0.073$	<b>0.795</b> <b><math>\pm 0.062</math></b>
NR-AR-LBD	<b>0.864</b> <b><math>\pm 0.056</math></b>	0.852 $\pm 0.050$	0.817 $\pm 0.057$	<b>0.858</b> <b><math>\pm 0.037</math></b>	0.816 $\pm 0.047$	0.844 $\pm 0.039$	<b>0.858</b> <b><math>\pm 0.027</math></b>	0.853 $\pm 0.034$
NR-AhR	<b>0.903</b> <b><math>\pm 0.028</math></b>	<b>0.900</b> <b><math>\pm 0.022</math></b>	0.854 $\pm 0.036$	0.896 $\pm 0.018$	0.869 $\pm 0.036$	0.890 $\pm 0.020$	<b>0.898</b> <b><math>\pm 0.020</math></b>	0.869 $\pm 0.023$
NR-Aromatase	0.827 $\pm 0.070$	0.802 $\pm 0.063$	0.742 $\pm 0.099$	0.824 $\pm 0.052$	<b>0.830</b> <b><math>\pm 0.050</math></b>	<b>0.845</b> <b><math>\pm 0.066</math></b>	<b>0.852</b> <b><math>\pm 0.053</math></b>	0.830 $\pm 0.065$
NR-ER	0.724 $\pm 0.019$	0.721 $\pm 0.022$	0.692 $\pm 0.018$	<b>0.734</b> <b><math>\pm 0.036</math></b>	<b>0.729</b> <b><math>\pm 0.025</math></b>	0.727 $\pm 0.043$	<b>0.733</b> <b><math>\pm 0.035</math></b>	0.712 $\pm 0.024$
NR-ER-LBD	<b>0.815</b> <b><math>\pm 0.051</math></b>	0.783 $\pm 0.055$	0.772 $\pm 0.019$	0.805 $\pm 0.024$	0.804 $\pm 0.030$	<b>0.810</b> <b><math>\pm 0.062</math></b>	<b>0.819</b> <b><math>\pm 0.036</math></b>	0.787 $\pm 0.038$
NR-PPAR-gamma	<b>0.839</b> <b><math>\pm 0.042</math></b>	0.793 $\pm 0.092$	0.756 $\pm 0.037$	<b>0.821</b> <b><math>\pm 0.105</math></b>	0.803 $\pm 0.064$	0.801 $\pm 0.104$	<b>0.825</b> <b><math>\pm 0.104</math></b>	0.783 $\pm 0.106$
SR-ARE	<b>0.818</b> <b><math>\pm 0.039</math></b>	<b>0.809</b> <b><math>\pm 0.039</math></b>	0.781 $\pm 0.046$	0.782 $\pm 0.040$	0.790 $\pm 0.049$	0.808 $\pm 0.028$	<b>0.826</b> <b><math>\pm 0.024</math></b>	0.777 $\pm 0.049$
SR-ATAD5	<b>0.857</b> <b><math>\pm 0.051</math></b>	0.828 $\pm 0.066$	0.738 $\pm 0.079$	<b>0.839</b> <b><math>\pm 0.037</math></b>	0.823 $\pm 0.041$	<b>0.841</b> <b><math>\pm 0.032</math></b>	0.837 $\pm 0.041$	0.811 $\pm 0.022$
SR-HSE	<b>0.793</b> <b><math>\pm 0.028</math></b>	0.764 $\pm 0.043$	0.731 $\pm 0.034$	<b>0.774</b> <b><math>\pm 0.036</math></b>	0.771 $\pm 0.036$	0.773 $\pm 0.049$	<b>0.786</b> <b><math>\pm 0.065</math></b>	0.750 $\pm 0.063$
SR-MMP	0.886 $\pm 0.019$	0.879 $\pm 0.026$	0.856 $\pm 0.027$	<b>0.888</b> <b><math>\pm 0.018</math></b>	0.886 $\pm 0.022$	<b>0.895</b> <b><math>\pm 0.017</math></b>	<b>0.909</b> <b><math>\pm 0.017</math></b>	0.865 $\pm 0.024$
SR-p53	<b>0.849</b> <b><math>\pm 0.033</math></b>	0.823 $\pm 0.057$	0.759 $\pm 0.034$	<b>0.840</b> <b><math>\pm 0.053</math></b>	0.813 $\pm 0.065$	0.833 $\pm 0.033$	<b>0.843</b> <b><math>\pm 0.054</math></b>	0.805 $\pm 0.036$
average	<b>0.830</b> <b><math>\pm 0.046</math></b>	0.811 $\pm 0.048$	0.771 $\pm 0.047$	0.821 $\pm 0.045$	0.811 $\pm 0.039$	<b>0.822</b> <b><math>\pm 0.045</math></b>	<b>0.831</b> <b><math>\pm 0.046</math></b>	0.803 $\pm 0.045$

- Mean and standard deviation on 5-fold cross validation.
- Top three results are **bolded**.
- MP graph is short for Message-passing graph.

**Figure A.1: Tukey's Test on Tox21**

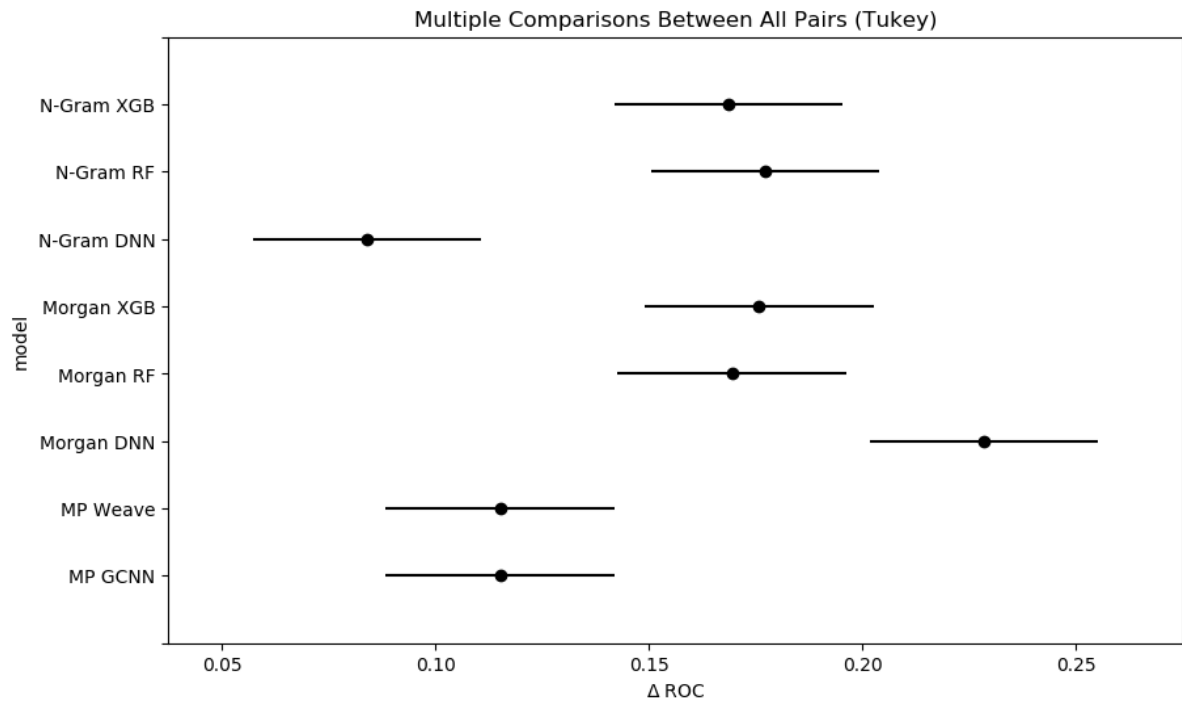


# Table 5: Generalization performance: Train and test gap on AUC[ROC]

	Morgan RF	Morgan XGB	Morgan DNN	MP graph GCNN	MP graph Weave	N-gram graph RF	N-gram graph XGB	N-gram graph DNN
NR-AR	0.213 $\pm 0.065$	0.209 $\pm 0.063$	0.243 $\pm 0.057$	<b>0.125</b> $\pm 0.085$	<b>0.115</b> $\pm 0.082$	0.198 $\pm 0.082$	0.210 $\pm 0.073$	<b>0.093</b> $\pm 0.072$
NR-AR-LBD	<b>0.136</b> $\pm 0.056$	0.144 $\pm 0.051$	0.180 $\pm 0.051$	<b>0.115</b> $\pm 0.040$	0.156 $\pm 0.036$	0.156 $\pm 0.039$	0.142 $\pm 0.027$	<b>0.112</b> $\pm 0.049$
NR-AhR	0.097 $\pm 0.028$	0.091 $\pm 0.023$	0.146 $\pm 0.036$	<b>0.054</b> $\pm 0.020$	<b>0.059</b> $\pm 0.045$	0.110 $\pm 0.020$	0.102 $\pm 0.020$	<b>0.048</b> $\pm 0.025$
NR-Aromatase	0.172 $\pm 0.070$	0.190 $\pm 0.064$	0.258 $\pm 0.099$	<b>0.112</b> $\pm 0.057$	<b>0.099</b> $\pm 0.064$	0.155 $\pm 0.066$	0.148 $\pm 0.053$	<b>0.066</b> $\pm 0.077$
NR-ER	0.274 $\pm 0.019$	0.238 $\pm 0.021$	0.307 $\pm 0.018$	<b>0.129</b> $\pm 0.033$	<b>0.129</b> $\pm 0.036$	0.273 $\pm 0.043$	0.267 $\pm 0.035$	<b>0.032</b> $\pm 0.026$
NR-ER-LBD	0.184 $\pm 0.051$	0.205 $\pm 0.056$	0.228 $\pm 0.019$	<b>0.134</b> $\pm 0.027$	<b>0.119</b> $\pm 0.061$	0.189 $\pm 0.062$	0.181 $\pm 0.036$	<b>0.105</b> $\pm 0.069$
NR-PPAR- gamma	0.161 $\pm 0.042$	0.199 $\pm 0.094$	0.244 $\pm 0.037$	<b>0.142</b> $\pm 0.110$	<b>0.147</b> $\pm 0.055$	0.197 $\pm 0.105$	0.175 $\pm 0.104$	<b>0.147</b> $\pm 0.129$
SR-ARE	0.181 $\pm 0.039$	0.166 $\pm 0.041$	0.219 $\pm 0.046$	<b>0.118</b> $\pm 0.053$	<b>0.099</b> $\pm 0.064$	0.192 $\pm 0.028$	0.174 $\pm 0.024$	<b>0.075</b> $\pm 0.080$
SR-ATAD5	0.143 $\pm 0.051$	0.167 $\pm 0.066$	0.262 $\pm 0.079$	<b>0.125</b> $\pm 0.042$	<b>0.129</b> $\pm 0.043$	0.159 $\pm 0.032$	0.163 $\pm 0.041$	<b>0.123</b> $\pm 0.032$
SR-HSE	0.206 $\pm 0.028$	0.222 $\pm 0.047$	0.269 $\pm 0.034$	<b>0.155</b> $\pm 0.052$	<b>0.155</b> $\pm 0.037$	0.225 $\pm 0.048$	0.214 $\pm 0.065$	<b>0.095</b> $\pm 0.074$
SR-MMP	0.114 $\pm 0.019$	0.109 $\pm 0.026$	0.144 $\pm 0.027$	<b>0.069</b> $\pm 0.029$	<b>0.063</b> $\pm 0.024$	0.105 $\pm 0.017$	0.091 $\pm 0.017$	<b>0.047</b> $\pm 0.037$
SR-p53	0.151 $\pm 0.033$	0.170 $\pm 0.058$	0.241 $\pm 0.034$	<b>0.107</b> $\pm 0.058$	<b>0.112</b> $\pm 0.069$	0.167 $\pm 0.033$	0.157 $\pm 0.054$	<b>0.064</b> $\pm 0.050$
average	0.169 $\pm 0.046$	0.176 $\pm 0.042$	0.229 $\pm 0.047$	<b>0.115</b> $\pm 0.027$	<b>0.115</b> $\pm 0.030$	0.177 $\pm 0.044$	0.169 $\pm 0.046$	<b>0.084</b> $\pm 0.033$

- Mean and standard deviation on 5-fold cross validation.
- Top three results are **bolded**.
- MP graph is short for Message-passing graph.

## Figure A.2: Tukey's Test on Generalization Gap



## Table A.1: Performance of Different Vector Embeddings

Original feature presented in baseline papers, and new feature presented in N-gram graph

	GCNN original feature	GCNN new feature	Weave original feature	Weave new feature
NR-AR	0.793 $\pm$ 0.068	0.802 $\pm$ 0.068	0.796 $\pm$ 0.057	0.792 $\pm$ 0.076
NR-AR-LBD	0.858 $\pm$ 0.037	0.852 $\pm$ 0.042	0.816 $\pm$ 0.047	0.830 $\pm$ 0.041
NR-AhR	0.896 $\pm$ 0.018	0.898 $\pm$ 0.018	0.869 $\pm$ 0.036	0.880 $\pm$ 0.034
NR-Aromatase	0.824 $\pm$ 0.052	0.819 $\pm$ 0.047	0.830 $\pm$ 0.050	0.834 $\pm$ 0.066
NR-ER	0.734 $\pm$ 0.036	0.735 $\pm$ 0.023	0.729 $\pm$ 0.025	0.734 $\pm$ 0.020
NR-ER-LBD	0.805 $\pm$ 0.024	0.803 $\pm$ 0.017	0.804 $\pm$ 0.030	0.809 $\pm$ 0.013
NR-PPAR-gamma	0.821 $\pm$ 0.105	0.799 $\pm$ 0.093	0.803 $\pm$ 0.064	0.801 $\pm$ 0.083
SR-ARE	0.782 $\pm$ 0.040	0.793 $\pm$ 0.046	0.790 $\pm$ 0.049	0.779 $\pm$ 0.022
SR-ATAD5	0.839 $\pm$ 0.037	0.838 $\pm$ 0.029	0.823 $\pm$ 0.041	0.800 $\pm$ 0.038
SR-HSE	0.774 $\pm$ 0.036	0.786 $\pm$ 0.030	0.771 $\pm$ 0.036	0.772 $\pm$ 0.038
SR-MMP	0.888 $\pm$ 0.018	0.887 $\pm$ 0.027	0.886 $\pm$ 0.022	0.885 $\pm$ 0.015
SR-p53	0.840 $\pm$ 0.053	0.836 $\pm$ 0.051	0.813 $\pm$ 0.065	0.813 $\pm$ 0.052
average	0.821 $\pm$ 0.045	0.821 $\pm$ 0.043	0.811 $\pm$ 0.039	0.811 $\pm$ 0.041

- We compare original GCNN and Weave and ones with feature vector in Eq (1)
- Here's the result for Tukey's Test:

## Table A.2: Tukey's Test on Models with Different Vector Embeddings

Group 1	Group 2	mean diff	reject
GCNN new feature	GCNN original feature	-0.0012	False
Weave new feature	Weave original feature	0.0008	False

- Null hypothesis is that both means are same. If reject=False, then we accept the null hypothesis: two means are same.
- This result shows that two vector embeddings (original paper and Eq (1) in the submitted version) have similar effects.

