

Deductive Machine Learning Models for Product Identification

Tianfan Jin, Qiyuan Zhao, Andrew B. Schofield, and Brett M. Savoie*

Davidson School of Chemical Engineering, Purdue University, West Lafayette, IN, 47906

E-mail: bsavoie@purdue.edu

1 Figures Referenced in the Main Text

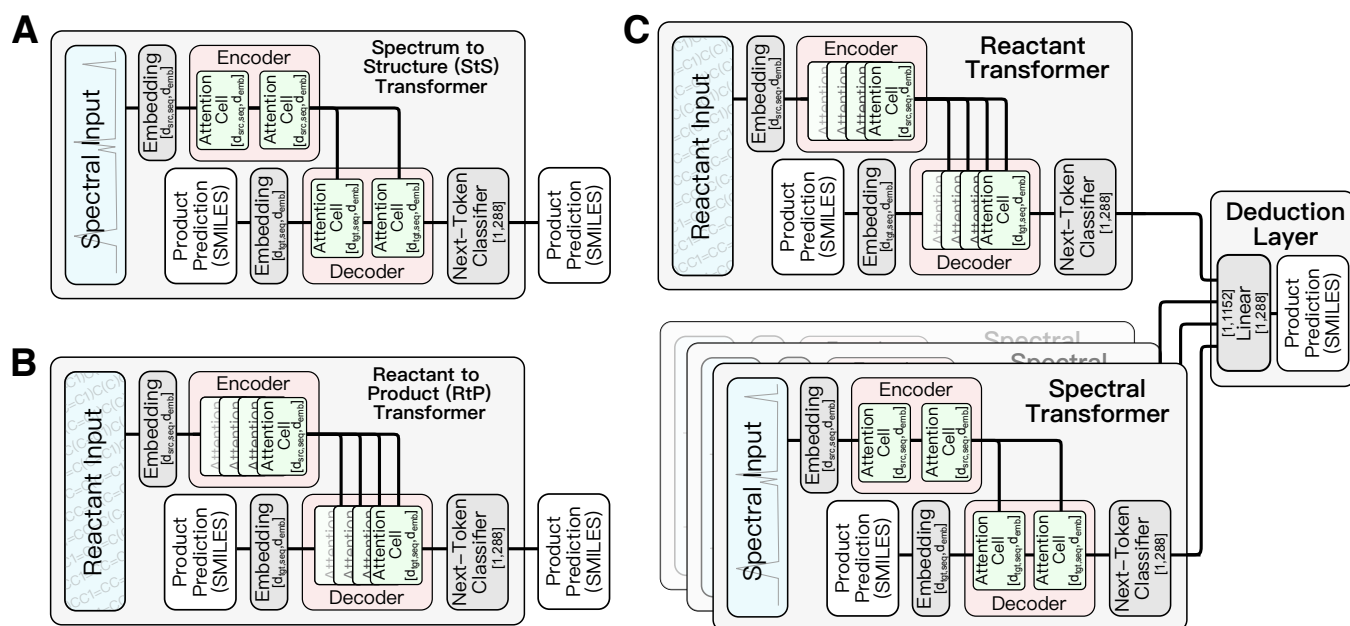


Figure S1: Overview of the different architectures developed in this work. (A) A sample spectrum to structure (StS) transformer. StS models with multiple spectral sources can be compared with the Reactant+Spectra models shown in (C), excluding the reactant transformer. (B) The reactant to product (RtP) model architecture. The use of four attention cells in the reactant encoder and decoder was found to significantly increase accuracy for this input. (C) The R+Spectra deductive architecture. The final linear layer projects to a token probability space using the token probabilities of the individual transformers. The input length $d_{\text{src}, \text{seq}}$ varies by input source (src), as described in the main text. The target (tgt) sequence length $d_{\text{tgt}, \text{seq}}$ is 64 and shared across models. The batch dimension of each model is omitted for clarity.

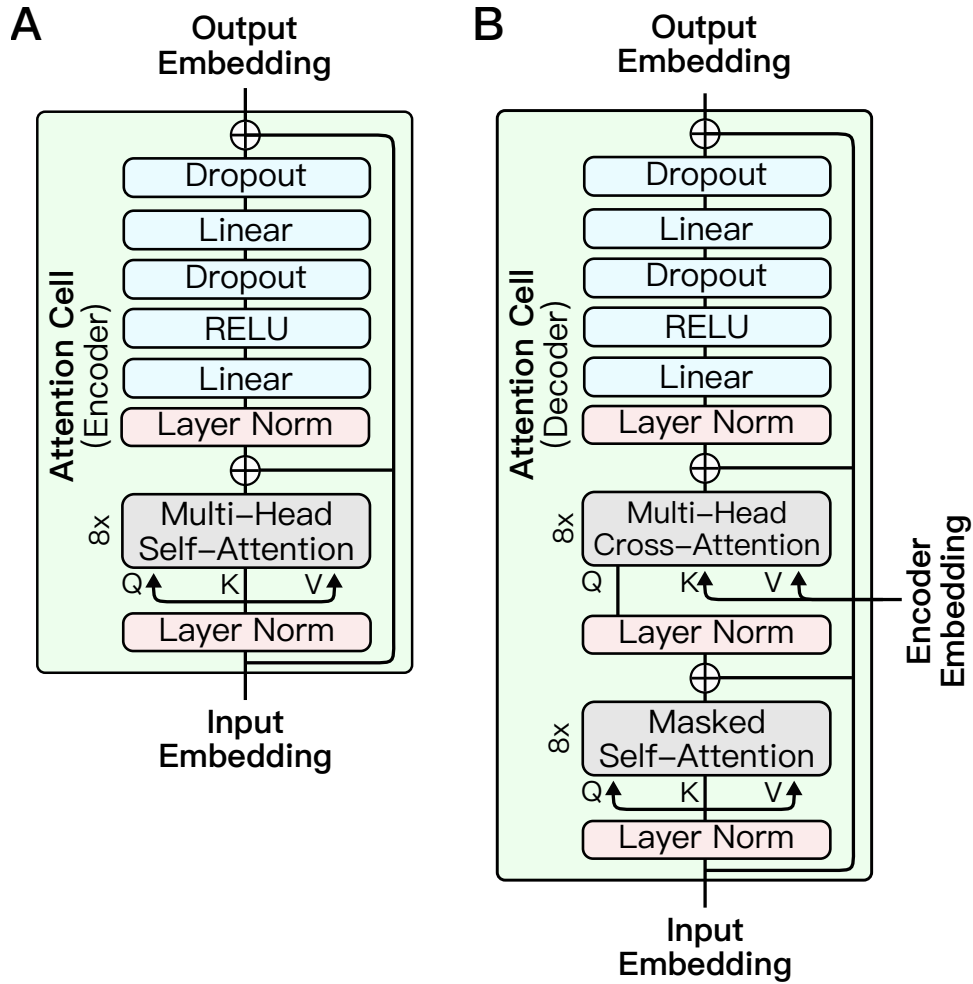


Figure S2: Overview of the attention cells used in the encoder and decoder of each transformer. (a) Encoder attention cell diagram. Each cell consists of typical multi-head attention, layer norm, residual connections, and feed-forward layers. (b) Decoder attention cell diagram. These cells are nearly-identical to the encoder attention cells, except that masking is used for the self-attention layer to limit information exchange to later tokens and there is a cross-attention layer inserted, whose key and value inputs are obtained as linear projections of the embedding dimension of the encoder output and the queries are obtained as linear projections of the embedding dimension of the output of the masked self-attention layer.

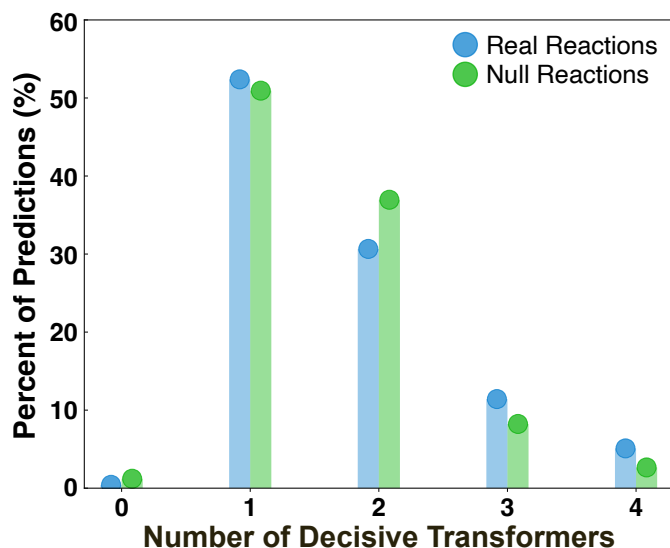


Figure S3: Results from decisiveness testing on the R+IR+NMR+MS model resolved by number of decisive transformers per product. A transformer was considered decisive for a product if it was decisive in decoding at least one token. Multiple transformers can be decisive for a given product, and cases with no decisive transformers mean that all tokens were governed by a consensus.

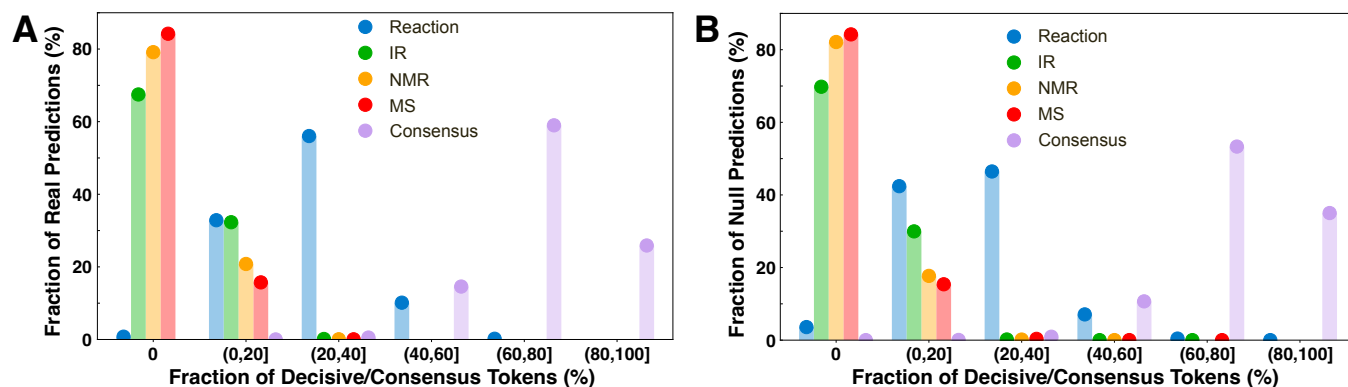


Figure S4: Results from decisiveness testing on the R+IR+NMR+MS model resolved on a per prediction and fraction of tokens basis. (A) Distribution of decisive tokens and consensus tokens across all testing samples corresponding to real product identification. (B) Distribution of decisive tokens and consensus tokens across all testing samples corresponding to starting material identification. A token was classified as a consensus prediction if no individual transformer was decisive in its prediction (i.e., at least two transformers supplied sufficient evidence to overrule the others in the top-1 prediction).

2 Illustrative Inference

We have provided an example to illustrate how the transformers act as deductive constraints upon each other through the recursive graph decomposition that occurs during inference (Fig. S5). The example, 1,4-dibromo-cyclohexane, is a molecule selected at random from the NIST chemistry WebBook and we have used the StS model for the illustration because the reactant-based models show less individual spectral decisiveness. The graph components are shaded based on the different information sources that were decisive in the inference.

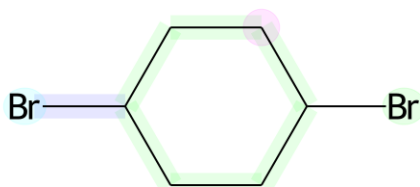


Figure S5: Inference of the MS+IR+NMR model on BrC1CCC(Br)CC1, a structure from the NIST WebBook. Each portion of the graph is shaded by transformer decisiveness: Green (MS decisive), Purple (NMR decisive), Cyan (NMR and MS decisive), Magenta (IR+NMR decisive).

The canonical SMILES for this species is BrC1CCC(Br)CC1. During inference, the starting “Br” is the first token that is decoded and this is almost entirely driven by the MS transformer. The other spectral transformers are not decisive for this token. However, the recursive manner in which inference occurs means that Br is now accepted by all transformers when decoding the rest of the molecular graph. Thus, the presence of the “Br” token is a constraint asserted by the MS that the rest of the transformers condition their inference on. The H-NMR and IR provide

other evidence for inference on the remainder of the graph (e.g., the joint IR/NMR decisiveness for the ring closure). Other notable behaviors include the order of the SMILES inference and the manner in which different spectral sources map to distinct molecular features. Although it is beyond the current scope to analyze in further detail here, we note that the inference patterns often map to typical expert heuristics associated with particular spectral sources. The inference order in this example matched the canonical SMILES, however more generally the model often prefers a non-canonical decoding if it has stronger confidence beginning on another portion of the graph. Similarly, the models often decode distinct SMILES associated with the same molecule in the top-n.

3 Learning Curves

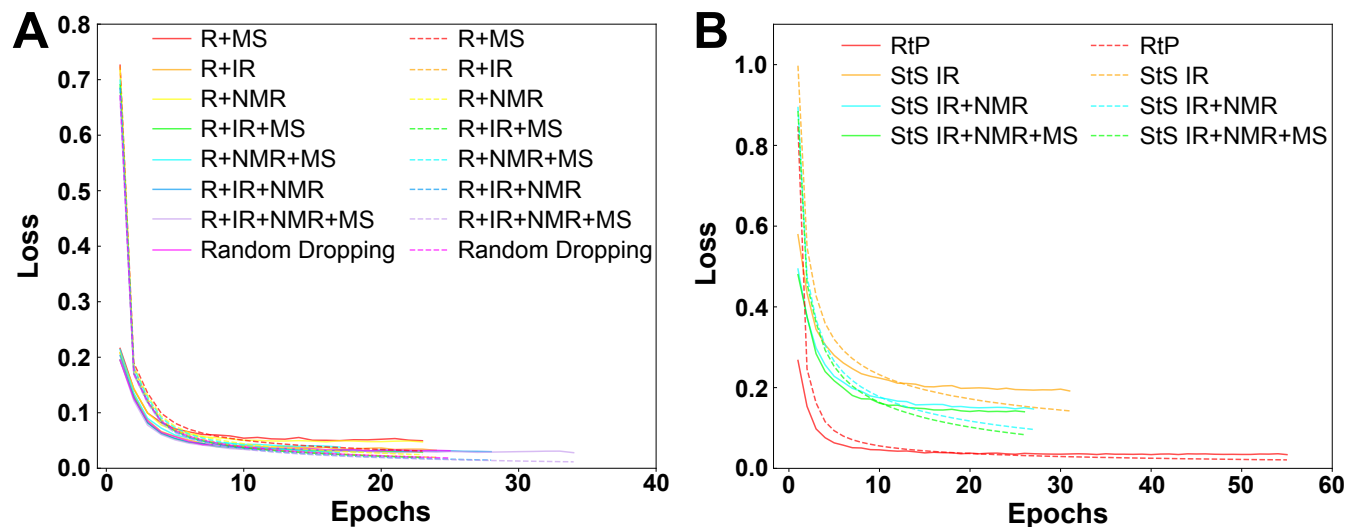


Figure S6: Learning curves for the models trained and tested in this work. Dotted lines correspond to the loss over training data and solid lines correspond to loss over validation data. (A) Composite of the reactant+spectra learning curves. (B) Composite of the StS and RtP learning curves.

4 MULTI Dataset

Table S1: Performance of the R+IR+NMR+MS model in identifying the products of reactions with unintended side-products. None of these examples were present in the USPTO dataset.

Reactants SMILES	Products SMILES	Prediction
<chem>CCC(C)Br.CC[O-]></chem>	<chem>CCOC(C)CC</chem>	Top-1
<chem>CCC(C)Br.CC[O-]></chem>	<chem>CC=CC</chem>	Top-1
<chem>O=C1CCCCC1.[N-]=[N+]=C></chem>	<chem>O=C1CCCCC1</chem>	Not in Top-5
<chem>O=C1CCCCC1.[N-]=[N+]=C></chem>	<chem>C1CCC2(CC1)CO2</chem>	Not in Top-5
<chem>CCOC1CCCCO1.CO></chem>	<chem>CCOC1CCCCO1</chem>	Top-1
<chem>CCOC1CCCCO1.CO></chem>	<chem>CCO</chem>	Top-1
<chem>O=C(O)c1ccccc1.CC=CC></chem>	<chem>O=C(O)c1ccccc1</chem>	Top-5
<chem>O=C(O)c1ccccc1.CC=CC></chem>	<chem>CC1OC1C</chem>	Not in Top-5
<chem>CS(=O)(=O)Cl.CN.CNC></chem>	<chem>CNS(C)(=O)=O</chem>	Top-5
<chem>CS(=O)(=O)Cl.CN.CNC></chem>	<chem>CN(C)S(C)(=O)=O</chem>	Top-1
<chem>Clc1ccccc1.CN.CNC></chem>	<chem>CNc1ccccc1</chem>	Top-1
<chem>Clc1ccccc1.CN.CNC></chem>	<chem>CN(C)c1ccccc1</chem>	Top-1
<chem>CN=C=O.CN.CNC></chem>	<chem>CNC(=O)N(C)C</chem>	Top-1
<chem>CN=C=O.CN.CNC></chem>	<chem>CNC(=O)NC</chem>	Top-5
<chem>CC=O.CCC(C)=O.CCNC.CN></chem>	<chem>CCC(C)NC</chem>	Top-1
<chem>CC=O.CCC(C)=O.CCNC.CN></chem>	<chem>CCN(C)CC</chem>	Top-1
<chem>CC=O.CCC(C)=O.CCNC.CN></chem>	<chem>CCC(C)N(C)CC</chem>	Not in Top-5
<chem>CC=O.CCC(C)=O.CCNC.CN></chem>	<chem>CCNC</chem>	Top-1
<chem>CC(=O)O.CN.CNC></chem>	<chem>CNC(C)=O</chem>	Top-1
<chem>CC(=O)O.CN.CNC></chem>	<chem>CC(=O)N(C)C</chem>	Top-1
<chem>COc1ccccc1I.N[Na]></chem>	<chem>COc1ccccc1N</chem>	Top-1
<chem>COc1ccccc1I.N[Na]></chem>	<chem>COc1cccc(N)c1</chem>	Not in Top-5
<chem>Clc1ccccc1.Nc1ccccc1.CN.CNC></chem>	<chem>CNc1ccccc1</chem>	Top-1
<chem>Clc1ccccc1.Nc1ccccc1.CN.CNC></chem>	<chem>CN(C)c1ccccc1</chem>	Top-1
<chem>Clc1ccccc1.Nc1ccccc1.CN.CNC></chem>	<chem>c1ccc(Nc2ccccc2)cc1</chem>	Top-1
<chem>CCC(=O)CC(C)=O.CNN></chem>	<chem>CCc1cc(C)nn1C</chem>	Not in Top-5
<chem>CCC(=O)CC(C)=O.CNN></chem>	<chem>CCc1cc(C)n(C)n1</chem>	Not in Top-5
<chem>CC(C)=O.CC=O.CCCL></chem>	<chem>CC=CC</chem>	Top-5
<chem>CC(C)=O.CC=O.CCCL></chem>	<chem>CC=C(C)C</chem>	Top-5
<chem>CC=CC.[O-][O+]=O></chem>	<chem>CC1OOC(C)O1</chem>	Not in Top-5
<chem>CC=CC.[O-][O+]=O>CCO</chem>	<chem>CCOC(C)OO</chem>	Not in Top-5
<chem>CC=CC.[O-][O+]=O>CCO</chem>	<chem>CC=O</chem>	Top-1
<chem>C=Pc1ccccc1.CC(C)=O></chem>	<chem>O=Pc1ccccc1</chem>	Not in Top-5
<chem>C=Pc1ccccc1.CC(C)=O></chem>	<chem>C=C(C)C</chem>	Top-1
<chem>CC=CC(O)COC(=O)c1ccccc1.CCOC(C)(OCC)OCC></chem>	<chem>CCOC(=O)CC(C)C=CCOC(=O)c1ccccc1</chem>	Not in top-5
<chem>CC=CC(O)COC(=O)c1ccccc1.CCOC(C)(OCC)OCC></chem>	<chem>CCO</chem>	Top-1
<chem>CC=CC></chem>	<chem>CC=CC</chem>	Top-5
<chem>CC=CCCC></chem>	<chem>CCCC=CCCC</chem>	Not in Top-5
<chem>CCc1ccccc1.BrBr></chem>	<chem>Br</chem>	Not in Top-5
<chem>CCc1ccccc1.BrBr></chem>	<chem>CC(Br)c1ccccc1</chem>	Top-1