

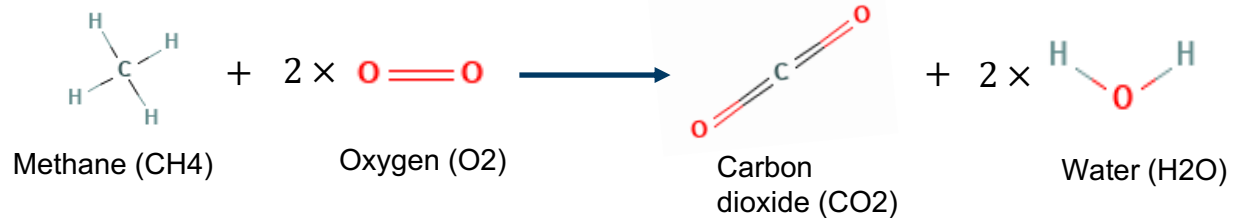
Representing chemical reactions

What is a chemical reaction?

- A chemical reaction is a process where substances (reactants/substrates) interact to form different substances (products)

- This is achieved by:

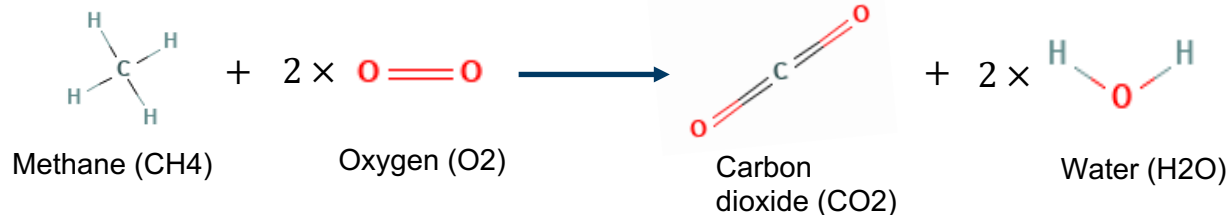
- Break bonds
- Forming bonds
- Rearranging bonds



- Law of conservation of mass: Matter cannot be created or destroyed in a chemical reaction
 - The number and type of atoms must remain constant throughout the reaction
- Different initiation or catalysis mechanisms:
 - Catalytic reactions
 - Enzyme-catalyzed reactions
 - Thermal reactions
 - ...

Using SMILES strings to represent reactions

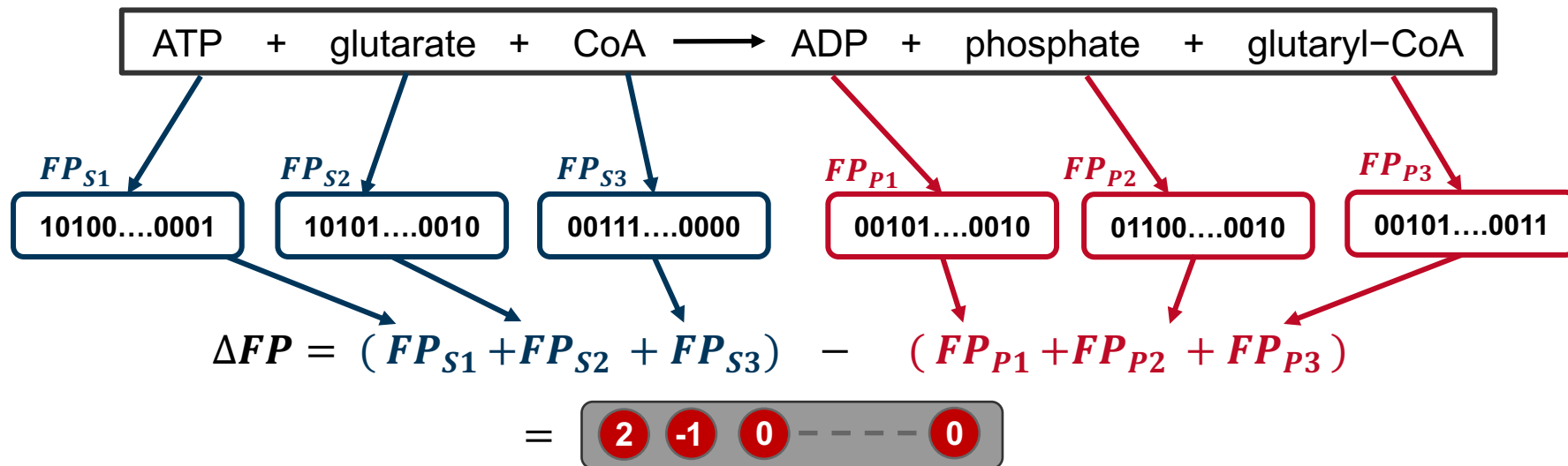
- We can use SMILES strings to represent the reactants and products
- We can combine SMILES strings to write down the whole reaction equations
 - “.” for separating the reactants and products
 - “>>” for separating reactants from products
- Example:



SMILES: C O=O O=C=O O

Reaction SMILES: C.O=O.O=O >> O=C=O.O.O

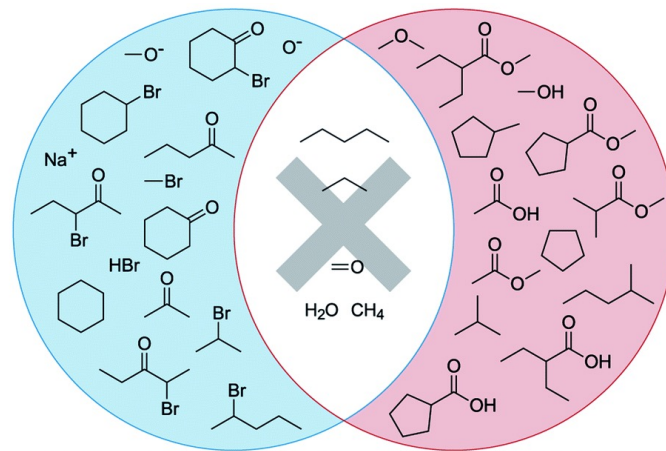
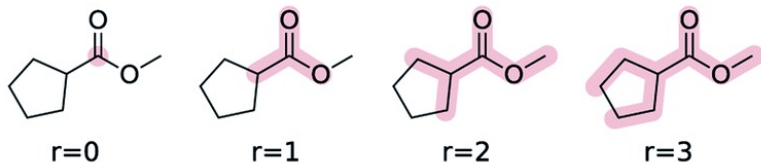
Traditional numerical representations – difference reaction fingerprints



Traditional numerical representations – DRFP



BrC1CCCCC1=O.C[O-].[Na+]>>COC(=O)C1CCCC1



{CC(=O)C, C[O-], C(C)(Br)C, ..., COC, CC1CCCC1, CC(C)C(=O)OC, CO}

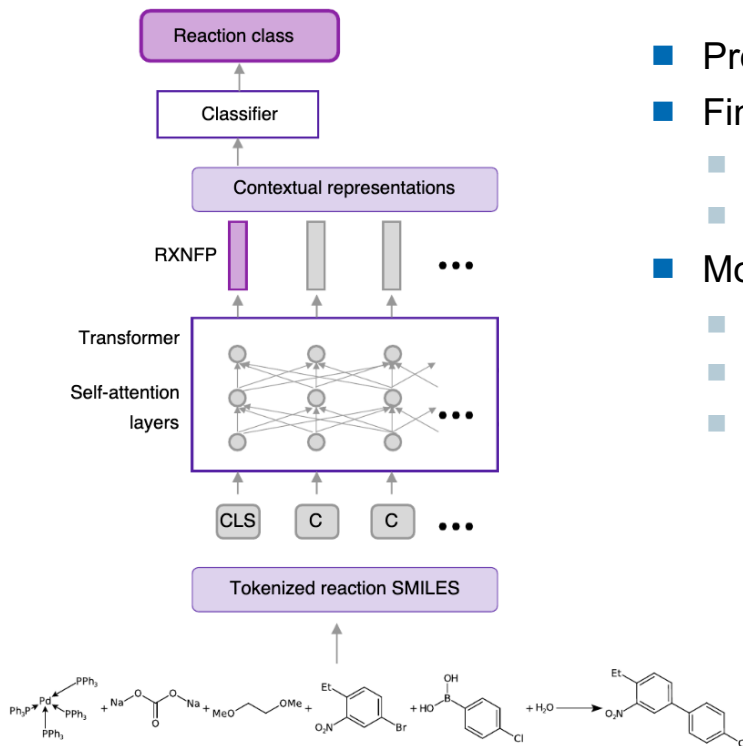
hashing: 32-bit hash function

[1440803970, 3834089465, 2582143990, ..., 322602819, 74077504]

folding: $x \bmod 2048$

0 1014 1154 1344 1859 2041 2047
[0, ..., 1, ..., 1, ..., 1, ..., 1, ..., 1, ..., 0]

Reaction Transformer Network (rxnfp)



- Pre-training: Masked Language Modelling (MLM)
- Fine-tuning: Predicting reaction classes
 - 2.6 million reaction equations
 - ~1000 reaction classes
- Model Architecture
 - Encoder-only model (BERT)
 - Hidden dimension $d = 256$
 - Number of layers encoder layer: