# **Supplementary Information for**

# Design of CO<sub>2</sub>-philic Molecular Units with Large Language Models

Konstantinos D. Vogiatzis\*

Department of Chemistry, University of Tennessee, Knoxville, Tennessee 37996-1600, United States

\*Corresponding Author: <u>kvogiatz@utk.edu</u>

# **Table of Contents**

Section	Page
S1. Prompts and Responses from ChatGPT-40	S-2
S2. Prompts and Responses from Llama-3	S-25
S3. Prompts and Responses from Gemini 2.0	S-27
S4. Computational Details	S-36
S5. List of Molecular Units and Computed Interaction Energies	S-38
S6. Additional Analysis of Best Performers	S-47
S7. Reproducability Experiments	S-49
S8. General Considerations and Best Practices	S-59
S9. Additional References	S-60

# S1. Prompts and Responses from ChatGPT-40

# January 25, 2025 (ChatGPT-40)

# **User (Initial Prompt):**

*An active field of research is the search of molecules with strong interaction with carbon dioxide* (CO2). These CO2-philic units can be introduced in materials such as metal-organic frameworks or polymers, and through enhanced interactions with CO2, they can be used for the selective separation and capture of CO2. For example, Vogiatzis et al. (ChemPhysChem 2009, 10, 374, https://doi.org/10.1002/cphc.200800583) estimated that the interaction energy between pyridine (SMILES: C1=CC=NC=C1) and CO2 is -4.46 kcal/mol, and between adenine (SMILES: C1=NC2=NC=NC(=C2N1)N) and CO2 is -5.90 kcal/mol. Similarly, Tian et al. (J. Phys. Chem. A 2015, 119, 3848, https://pubs.acs.org/doi/10.1021/acs.jpca.5b01892) found that the interaction energy between triazole-based molecules that can go up to -7.00 kcal/mol. More recently, Townsend et al. (Nat. Commun. 2020, 11, 3230, https://www.nature.com/articles/s41467-020-17035-5) used machine learning and discover new organic molecules with enhanced CO2philicity. For example, the molecule with SMILES string C1(=C(N=C([N]1[H])O[H])N([H])[H])N([H])[H] has a predicted interaction energy with CO2 of -7.41 kcal/mol. However, for enhance performance, it has been hypothesized that the interaction energy between these organic molecules and CO2 should be lower, at about -10 kcal/mol.

I would like to ask you the following two questions:

- 1. Can you compile additional data from the literature with similar CO2-philic groups. You can either provide their SMILES string or their molecular formula. Listing the corresponding references would be much appreciated.
- 2. Based on the data that you have found and read, can you predict molecules with CO2 interaction that is lower than -7.5 kcal/mol and closer to -10 kcal/mol?

Thanks :)

# Response #1 (GPT-40):

The search for molecules with strong interactions with carbon dioxide  $(CO_2)$  is a significant area of research, particularly for applications in  $CO_2$  capture and separation. Identifying  $CO_2$ -philic groups that can be incorporated into materials like metal-organic frameworks (MOFs) or polymers is crucial for enhancing their performance.

# 1. Compilation of CO<sub>2</sub>-Philic Groups:

Several studies have identified functional groups and molecules with notable affinity for CO<sub>2</sub>:

• **Poly(ethylene oxide) (PEO) Oligomers:** Ab initio calculations have shown that PEO oligomers with more than three repeating units exhibit significant CO<sub>2</sub>-binding affinity.

This finding aligns with the high performance of PEO-based materials in  $\rm CO_2/N_2$  separation processes.

# pubs.acs.org

• Acylamide-Functionalized Ligands: A flexible C<sub>3</sub>-symmetric hexacarboxylate ligand with acylamide groups, when incorporated into an rht-type MOF, demonstrated increased CO<sub>2</sub> uptake and enhanced heat of adsorption compared to analogous MOFs with alkyne groups.

# pubs.acs.org

• Amine Groups: It is well-established that CO<sub>2</sub> molecules have a high affinity for amine groups, leading to the formation of carbamates. This interaction is utilized in various CO<sub>2</sub> capture technologies.

# Nature

• **Carbonyl and Ether Groups:** Research indicates that the binding energy between CO<sub>2</sub> and carbonyl oxygen is comparable to that between CO<sub>2</sub> and ether oxygen. This suggests that polymers containing these groups could interact favorably with CO<sub>2</sub>.

# **D-Scholarship**

# 2. Predicting Molecules with Enhanced CO<sub>2</sub> Interaction:

Building upon the data from these studies, molecules containing multiple CO<sub>2</sub>-philic groups may exhibit stronger interactions with CO<sub>2</sub>. For instance, combining amine and carbonyl functionalities within a single molecule could potentially enhance CO<sub>2</sub> binding affinity. Additionally, increasing the number of repeating units in PEO oligomers beyond three may further strengthen CO<sub>2</sub> interactions.

While specific interaction energies are not provided in these studies, the presence of multiple CO<sub>2</sub>-philic groups within a molecule is likely to result in interaction energies lower (more negative) than -7.5 kcal/mol, potentially approaching the -10 kcal/mol threshold hypothesized for enhanced performance.

Further computational and experimental studies would be necessary to quantify the interaction energies of such designed molecules and confirm their efficacy in CO<sub>2</sub> capture applications.

<u>User:</u> Thank you. I would like to perform such computations to discover new molecules with stronger CO2 affinity. However, since the chemical space is vast, I would like from you to utilize the knowledge that you shared in your response and write down molecules (in SMILES format) that can have interaction energy closer to -10 kcal/mol.

# Response #2:

Based on the CO<sub>2</sub>-philic groups identified from the literature, here are some potential molecules in SMILES format that might exhibit stronger interactions with CO<sub>2</sub>, potentially approaching the -10 kcal/mol threshold:

- 1. Tris(2-aminoethyl)amine (TREN) with carbonyl groups:
  - N(CCN(CCN)C=O)C=O
- 2. Polyethylene oxide (PEO) oligomer with amine end-groups: • NCCOCCOCCOCCOC
- 3. Hexa-acylamide functionalized benzene core:
  - C1(C(=0)NC2=CC=C(C=C2)C(=0)NCC(=0)N
- 4. Triazole-functionalized polyethylene glycol (PEG): o N1C=CN=C1CCOCCOCCOCCN1C=CN=C1
- 5. Ethylenediamine-tetraacetic acid (EDTA) analogue with carbonyl groups: o C(CN(CC(=O)O)CC(=O)O)(CN(CC(=O)O)CC(=O)O)N
- 6. Amine-functionalized cyclodextrin:
  - C1OC2(CO1)OC(CO)C(O2)N

These molecules incorporate various functional groups known for their CO<sub>2</sub>-philic nature, such as amines, carbonyls, triazoles, and ether groups. Computational methods such as density functional theory (DFT) or molecular dynamics simulations could be used to estimate their CO<sub>2</sub> interaction energies and refine molecular designs.

<u>User:</u> Thank you. However, these are known molecules. In the field of molecular discovery, we often explore not-yet-synthesized molecules with the aim to have increased performance for a specific chemical application. Could you please write down not-yet-synthesized molecules (in SMILES format) that can have interaction energy closer to -10 kcal/mol.

# Response #3:

Got it! Here are some hypothetical, not-yet-synthesized molecules that incorporate multiple CO<sub>2</sub>-philic functional groups (such as amines, carbonyls, and ethers) strategically placed to enhance interactions with CO<sub>2</sub>. These molecules are designed to potentially exhibit interaction energies closer to -10 kcal/mol:

- 1. Hypothetical amide-ether hybrid framework:
  - N(CCOCCOCC(=0)NC(=0)OCCOCCN)C(=0)NC(=0)O
- 2. Branched polyamine-triazole system: • N1C=CN=C1CC(CN(CCOCCOCCN)CCN1C=CN=C1)C(=O)O
- 3. Cyclic multi-ether-multi-amine scaffold:
  - C1(CCOCCOCCNCCOCCNCCOCCO)N1
- 4. Extended polyethylene glycol (PEG) backbone with triazole and amide terminations:
  - N1C=CN=C1CCOCCOCCOCC(=0)N(CCN1C=CN=C1)C(=0)O
- 5. Hyperbranched amine-ether-carbonyl polymer fragment: • N(CCOCCN(CC(=O)NC(=O)OCCOCCN)C(=O)O)C(=O)NC(=O)O
- 6. **Multifunctional cyclized triazole-amine-carbonyl structure:** • C1(=CC(=O)N(CCN1C=CN=C1)C(=O)N(CCOCCOCCN)C(=O)O)N1C=CN=C1
- 7. Dendrimeric amine-triazole-carbonyl construct:
  - N(CCN(CC(=0)NC(=0)N1C=CN=C1)C(=0)0)C(=0)0)C(=0)0

These structures strategically combine functionalities known to enhance CO<sub>2</sub> affinity while incorporating novel connectivity patterns and branching to increase surface area and binding sites. Computational screening could help assess their feasibility and optimize their design for better CO<sub>2</sub> capture performance.

<u>User:</u> Thank you, this is useful. However, most of your structures are amine-based structures, while nitrogen-containing heterocyclic structures have typically stronger noncovalent interactions than amines. Could you please write down not-yet-synthesized molecules (in SMILES format) that can have interaction energy closer to -10 kcal/mol and are not based on amines.

# Response #4:

Thanks for the clarification! Here are some hypothetical nitrogen-containing heterocyclic structures (excluding amine-based functionalities) that could potentially exhibit CO<sub>2</sub> interaction energies closer to -10 kcal/mol:

- 1. Fused triazole-pyrimidine system with hydroxyl substitutions:
  - C1=C(N=CN=C10)C2=NC=NC(=C20)C3=NC=NC=C30
- 2. Tetrazole-pyrazine hybrid with multiple hydroxyl groups: o C1=NC(=NC=N1)C2=NC=CN=C2OCCO
- 3. Benzotriazole-furan fused system:
  - C1=CC2=NN=CN2C=C1OCCOCCO
- 4. **Polycyclic oxadiazole-triazine structure:** • C1=NC(=NO1)C2=NC=NC=N2C3=NC=NO3
- 5. Functionalized quinoxaline-tetrazole framework: • C1=NC2=CC=NC=C2N=C1C3=NN=NN3C(=O)O
- 6. Imidazole-pyrazole-thiadiazole hybrid: • C1=NC=CN=C1C2=CN=NN2C3=NC=CS3
- 7. **Triazine-benzimidazole extended network:** • C1=NC(=NC(=N1)C2=NC3=CC=CC=C3N2)OCCO
- 8. Diazaphenanthrene-triazole fusion:
  - C1=NC2=CC=NC=C2C3=NN=CN=C31

These structures incorporate highly polar nitrogen-containing heterocycles such as triazoles, tetrazoles, and triazines, which have shown strong noncovalent interactions with CO<sub>2</sub> in previous studies. They also include additional polar functionalities like hydroxyl and ether groups to further enhance CO<sub>2</sub> binding affinity.

These molecules could be excellent candidates for computational screening to predict their interaction energies and assess their suitability for CO<sub>2</sub> capture applications.

<u>User:</u> Thank you, this is great! We also know that adenine-based structures have higher CO2philicity than triazoles or triazines. Could you please write down not-yet-synthesized molecules (in SMILES format) that can have interaction energy closer to -10 kcal/mol and utilize nucleases such as adenine?

# Response #5:

You're welcome! Here are some hypothetical adenine-based molecules designed to enhance CO<sub>2</sub> affinity, potentially achieving interaction energies closer to -10 kcal/mol:

- 1. **Hydroxy-functionalized adenine dimer with ether linkage:** o C1=NC2=C(N=CN2C=N1)OCCOCC3=NC4=C(N=CN4C=N3)O
- 2. Adenine-tetrazole hybrid with carboxyl substitutions:
  - C1=NC2=C(N=CN2C=N1)C3=NN=NN3C(=O)O
- 3. Fused adenine-pyrazole-quinoxaline system: • C1=NC2=C(N=CN2C=N1)C3=NC4=CC=NC=C4N3
- 4. Adenine-benzimidazole scaffold with ether and hydroxyl groups:
  - C1=NC2=C(N=CN2C=N1)N3C=NC4=CC=CC=C43OCCO
- 5. Branched adenine-polyether framework: o C1=NC2=C(N=CN2C=N1)OCCOCCOCCOC
- 6. Trimeric adenine with hydroxyl and carbonyl groups: • C1=NC2=C(N=CN2C=N1)C3=NC4=C(N=CN4C=N3)C5=NC6=C(N=CN6C=N5)OCC(=O)O
- 7. Adenine-thiadiazole hybrid for increased sulfur interactions: • C1=NC2=C(N=CN2C=N1)C3=NC=NS3
- 8. Extended adenine scaffold with multiple ether groups:
  - C1=NC2=C(N=CN2C=N1)OCCOCCOCCOCCO

These structures strategically incorporate adenine nucleobases to leverage their strong CO<sub>2</sub> affinity while introducing additional polar functional groups (hydroxyl, carbonyl, ether) to further enhance binding interactions. Such molecules could be promising candidates for computational screening to optimize their CO<sub>2</sub> capture performance.

<u>User:</u> Thank you very much. I will perform DFT calculations with these molecules and hopefully I will get stronger interaction energy with CO2.

# Response #6

That sounds like a great plan! DFT calculations should provide valuable insights into the interaction energies and help fine-tune the molecular structures for even better CO<sub>2</sub> affinity. Let me know if you need any help with setting up calculations, choosing functionals/basis sets, or analyzing the results.

Good luck with your research—hope you discover some highly CO<sub>2</sub>-philic molecules!

<u>User:</u> You gave me two SMILES strings that have mistakes and I cannot convert them to xyz coordinates. Could you please fix them? Those are:

*C1=NC(=NO1)C2=NC=NC=N2C3=NC=NO3* 

*C1=NC2=C(N=CN2C=N1)N3C=NC4=CC=CC=C43OCCO* 

# Response #7

Thanks for pointing that out! Let me check and correct the SMILES strings for you.

# **Fixed SMILES strings:**

- Original: C1=NC(=NO1)C2=NC=NC=N2C3=NC=NO3
   Correction: C1=NC(=NO1)C2=NC=NC=N2C3=NC=N[O]3
   (The issue might be with the placement of oxygen in the third ring; using [O] ensures
   correct atom placement.)
- Original: C1=NC2=C(N=CN2C=N1)N3C=NC4=CC=CC=C43OCCO Correction: C1=NC2=C(N=CN2C=N1)N3C=NC4=CC=CC=C4-3OCCO (The issue could be a missing ring closure; adding a dash - ensures proper connectivity.)

Try using these corrected versions in your conversion script. Let me know if you encounter any further issues!

<u>User:</u> Thank you. These two SMILES strings have similar issues. Could you please fix those as well?

# C1(CCOCCOCCNCCOCCOCO)N1

C1(=CC(=O)N(CCN1C=CN=C1)C(=O)N(CCOCCOCCN)C(=O)O)N1C=CN=C1

# Response #8:

Sure! Here are the corrected versions of your SMILES strings:

- Original: C1(CCOCCOCCNCCOCCO)N1
   Correction: C1(CCOCCOCCNCCOCCO)N1 → Likely issue: improper ring closure.
   Fixed: N1(C(CCOCCOCCNCCOCCO)C1)
   (Ensures proper ring formation around nitrogen.)
- Original: C1(=CC(=O)N(CCN1C=CN=C1)C(=O)N(CCOCCOCCN)C(=O)O)N1C=CN=C1 Correction: Possible issue with nested ring numbering and closure. Fixed: C1=CC(=O)N(CCN1C=CN=C1)C(=O)N(CCOCCOCCN)C(=O)O.N1C=CN=C1 (Splitting into separate fragments might fix structural issues.)

Try using these corrected versions and let me know if you need further adjustments!

# February 5, 2025 (ChatGPT-4o and reasoning)

<u>User:</u> A few days ago, we discussed how we can design new molecules for enhancing CO2 capture. Our target is to find molecules that have interaction energy lower than -10 kcal/mol. You provided 29 different molecules (as SMILES strings), in groups of 4. You can check the history of our conversation for more details. I performed accurate quantum chemical calculations and I found that the fourth group (adenine-based organic molecules) resulted to lower interaction energies, with an average of -5.9 kcal/mol, and the

C1=NC2=C(N=CN2C=N1)C3=NC4=C(N=CN4C=N3)C5=NC6=C(N=CN6C=N5)OCC(=O)Omolecules was the best performer (-7.0 kcal/mol). However, we are still far from our target (lower than -10 kcal/mol). I would like to use the knowledge that you have acquired and provide a list of 20 molecules (as SMILES strings) with potential lower interaction energy than -10 kcal/mol. I also add a list of design parameters based on my intuition that you can consider in your molecular design task:

- nitrogen atoms that are part of heteroaromatic molecular groups lower the interaction energy
- If we have two of such nitrogen containing heteroaromatic groups where both nitrogen atoms interact with CO2, this lowers the interaction energy, e.g.
   C1=NC2=C(N=CN2C=N1)C3=NC4=C(N=CN4C=N3)C5=NC6=C(N=CN6C=N5)OC C(=O)O
- Molecular groups with hydrogen atoms (e.g. OH, NH2) that are close to the nitrogen that interacts with CO2, are also helping and lower the interaction energy
- Electron donating groups that are part of the heteroaromatic molecular groups will increase the electron density of the nitrogen that interacts with CO2, and this should further lower the interaction energy
- Lengthy molecular chains such as N(CCN(CCN)C=O)C=O do not help since they tend to block the nitrogen that interacts with CO2

# Response #9:

*(User utilized the option "reason", response in 27 seconds)* 

# Mapping out molecule design

I'm gathering ideas to design molecules with lower CO2 capture energy. Adenine-based molecules show promise, but we still need molecules with interaction energy below -10 kcal/mol. Adding nitrogen and hydrogen functionalities might help.

#### Mapping out potential candidates

I'm considering heteroaromatic structures with nitrogen and electron-donating groups near nitrogen centers. Long chains may hinder CO2 interaction, so shorter, more compact molecules are preferable.

# Focusing on design

I'm narrowing down heteroaromatic structures with nitrogen and electron-donating groups. Avoiding long chains, I'm leaning towards bipyrimidine derivatives for their potential low CO2 interaction energy.

# **Crafting bipyrimidine derivatives**

I'm working on creating 20 bipyrimidine derivatives with amino and hydroxy substituents. Initially, I explored adenine, but now I'm focusing on connected heteroaromatics with NH2 and OH groups near nitrogen.

# **Constructing bipyrimidine derivatives**

I'm piecing together 20 bipyrimidine derivatives with amino and hydroxy substituents near interacting nitrogen atoms, aiming for optimal CO2 interaction energy.

# Generating distinct molecules

I'm preparing to create 20 unique SMILES strings, spanning bipyrimidine, bipyridine, bipyrazine, triazine, pyrimidopyrimidine, and melamine derivatives, ensuring both nitrogen atoms remain accessible.

# **Crafting molecule structures**

I'm formulating bipyrimidine and bipyrazine derivatives with amino and hydroxyl groups, ensuring smooth interaction with CO2 and maintaining accessible nitrogen atoms. Ready to generate distinct SMILES strings.

# Generating molecular candidates

I'm crafting varied derivatives like bipyrimidine, bipyrazine, triazine, and bipodal melamine, each with distinct substituents to ensure accessible nitrogen atoms and optimal CO2 interaction energies.

# Formulating potential derivatives

I'm piecing together bipyrimidine derivatives with varied substituents, considering simpler structures for easier nitrogen atom accessibility, and aiming for optimal CO2 interaction energies.

# **Evaluating molecular candidates**

I'm considering fused ring systems like purine derivatives with NH2 groups, and exploring molecules with three heteroaromatic rings, though complexity and substituent positions are challenges.

# **Mapping distinct SMILES creations**

I'm listing various SMILES for bipyrimidine, bipyrazine, bipodal melamine, and pyrimidinepyridine derivatives. Each molecule combines different substituents, aiming for unique and plausible structures.

# Expanding molecule designs

Pulling together bipyrimidine dimethylamino, bipyrazine derivatives with varied substituents, and purine and triazine derivatives, each optimized for CO2 interaction and nitrogen atom accessibility.

# Verifying SMILES strings

I'm meticulously ensuring each SMILES string for derivatives with triazine, s-triazine, pyrimidine, and fused heteroaromatic systems is valid and optimized for CO2 interaction and nitrogen atom accessibility.

# **Crafting final SMILES**

I'm developing the final 20 SMILES strings for hypothetical derivatives, ensuring each meets the CO2 interaction and nitrogen accessibility parameters, while maintaining SMILES validity.

# **Ensuring completeness**

I'm finalizing the list, ensuring each SMILES string is accurate and adheres to CO2 interaction and nitrogen accessibility criteria. This thorough review confirms the list's robustness.

# **Ensuring accuracy**

OK, let me see. I'm making sure each SMILES string is correct and ready for the final list of 20 candidates. This ensures they align with CO2 interaction energy and nitrogen accessibility.

# **Ensuring accuracy**

Verifying SMILES strings for derivatives like bipyrimidine, bipyrazine, bipodal melamine, and fused purine derivatives. Ensuring alignment with CO2 interaction and nitrogen accessibility parameters.

# **Refining the list**

I'm ensuring the completeness and CO2 interaction energy of the remaining SMILES strings, focusing on chemical validity and specificity to boost CO2 interaction.

# **Confirming chemical structures**

Ensuring each SMILES string is chemically valid and aligned with CO2 interaction criteria. The list includes bipyrimidine, bipyrazine, bipodal melamine, and fused purine derivatives.

Below is a list of 20 hypothetical candidate molecules (given as SMILES strings) that combine multiple nitrogen-containing heteroaromatic rings decorated with nearby  $-NH_2$  or -OH (or even  $-NCH_3$ ) groups. These designs follow your guidelines: they feature (a) heteroaromatic rings in which ring N-atoms are available for CO<sub>2</sub> interaction; (b) two such rings in many cases (so that both "active" N-sites can interact); (c) substituents ( $-NH_2$ , -OH, or  $-NCH_3$ ) in positions close to the ring N-atoms to boost the local electron density; and (d) they avoid long flexible chains that might sterically hinder the binding. (Keep in mind that these are "in-silico" candidates suggested on the basis of design intuition and must be validated by further calculations/experiments.)

# 1. Bipyrimidine diamine

n1c(N)nc(c1)-c2nc(N)nc(n2)

# 2. Bipyrimidine dihydroxy

n1c(O)nc(c1)-c2nc(O)nc(n2)

# 3. Bipyrimidine mixed –NH<sub>2</sub>/–OH

n1c(N)nc(c1)-c2nc(O)nc(n2)

# 4. Bipyrazine diamine

c1c(nccn1N)-c2c(nccn2N)

# 5. Bipyrazine dihydroxy

clc(nccn1O)-c2c(nccn2O)

# 6. Bipodal melamine (s-triazine) diamine

Nc1nc(N)nc(N)c1-c2nc(N)nc(N)c2

# 7. Bipodal melamine with mixed -NH<sub>2</sub>/-OH

Nc1nc(N)nc(N)c1-c2nc(N)nc(O)c2

# 8. Pyrimidine–pyridine diamine

n1c(N)nc(c1)-c2ncccn2

# 9. Pyrimidine–pyridine dihydroxy

n1c(O)nc(c1)-c2ncccn2

# 10. Pyrimidine-pyridine mixed (-NH<sub>2</sub> on one ring, -OH on the other)

n1c(N)nc(c1)-c2ncccn2O

# 11. Bipyrimidine dimethylamino derivative

n1c(NC)nc(c1)-c2nc(NC)nc(n2)

# 12. Bipyrimidine with an extra -NH2 on one ring

n1c(N)nc(c1)-c2nc(N)nc(c2)N

# 13. Fused purine derivative with an -NH2 substituent

clncnc2ncn(N)c2n1

# 14. Fused purine derivative with an -OH substituent

c1ncnc2ncn(O)c2n1

# 15. Alternative fused purine with -NH2 at a different position

c1ncnc2n(cn1)nc(N)c2

# 16. Bipodal s-triazine (melamine-like) diamine

Nc1ncncn1-c2ncncn2

# 17. Bipodal s-triazine dihydroxy

Oc1ncncn1-c2ncncn2

# 18. Pyrimidine-s-triazine mixed (-NH2 on the pyrimidine ring)

n1c(N)nc(c1)-c2ncncn2

# 19. Pyrimidine-s-triazine mixed (-OH on the pyrimidine ring, -NH2 on the triazine)

n1c(O)nc(c1)-c2ncncn2N

# 20. Fused bicyclic heterocycle with -NH2 and -OH substituents

n1c(N)nc2c1ncnc2O

Each of these SMILES strings represents a compact, rigid heteroaromatic system that (by design) should expose electron-rich nitrogen sites for favorable  $CO_2$  binding. Although these molecules are proposed to potentially achieve interaction energies below -10 kcal/mol, further quantum-chemical validation is recommended.

# User:

In the following SMILES, you have considered molecules that have a nitrogen atom that has four bonds. There are not acceptable chemical structures. clc(nccn1N)c2c(nccn2N) clc(nccn1O)c2c(nccn2O) nlc(N)nc(c1)-c2ncccn2O Nc1ncncn1-c2ncncn2 Oc1ncncn1c2ncncn2 nlc(O)nc(c1)c2ncncn2N

Instead of placing the amino group (NH2) or the hydroxo group (OH) directly bonded to N of the aromatic ring, could you please generate 20 more molecules (as SMILES strings) that have these groups on the carbon that is next to the N of the aromatic ring that can interact with CO2?

# Response #10:

(User utilized the option "reason", response in 43 seconds)

# **Designing molecules**

I'm refining the molecule design to place substituents (-NH2 or -OH) on the carbon atoms adjacent to the nitrogen atoms in heteroaromatic rings, as this configuration lowers interaction energy with CO2.

# **Crafting bipyrimidine**

I'm working on 2-bipyrimidine derivatives, focusing on substituents (like 2-aminopyrimidine or 2-hydroxy) attached to carbons next to nitrogen. This approach should effectively lower interaction energy with CO2.

# **Pinpointing atom placements**

I'm thinking through pyrimidine atom numbering, ensuring substituents are on carbons next to nitrogen atoms and moving systematically to design the molecules.

# Mapping the structure

I'm piecing together the detailed structure of 2-aminopyridine, ensuring the amino group is on the carbon next to the nitrogen, while carefully considering the aromatic ring and nitrogen interactions.

# **Aligning positioning**

I'm confirming the amino group is on the carbon adjacent to nitrogen in 2-aminopyridine, clarifying its exact spatial relationship within the molecule.

# **Clarifying substituent placement**

I'm aligning the amino group's position on 2-aminopyridine's carbon atom next to nitrogen, ensuring precise placement in the SMILES notation, and confirming this arrangement.

# **Ensuring precision**

Adjusting my approach to ensure substituent placement next to nitrogen in SMILES notation for bipyridines or bipyrimidines, aiming for accurate molecular structure representation.

# Mapping molecule structures

I'm planning to synthesize 20 molecules, focusing on bipyridines or bipyrimidines, with amino groups adjacent to ring nitrogen, ensuring precise molecular representation in SMILES notation.

# **Ensuring precise bond placement**

Confirming the amino group is adjacent to nitrogen in bipyridine and bipyrimidine structures, aiming for accurate molecular representation in SMILES notation by focusing on precise bonding patterns.

# Mapping molecular design

I'm thinking through bipyridine structures with amino groups near nitrogen, avoiding long linkers to keep functional groups close for optimal CO2 interaction. Progress is steady and purposeful.

# Enhancing molecular bonding

I'm detailing bipyridine and bipyrimidine structures with precise amino and hydroxyl group placements, ensuring functional groups are near nitrogen atoms for optimal CO2 interaction, avoiding long linkers.

# **Determining carbon placement**

I'm thinking through bipyrimidine structures to strategically place amino and hydroxyl groups near interacting nitrogen. Ensuring clarity in carbon attachment and bond linkage to optimize molecular interactions.

# Crafting the list

I'm assembling 20 molecules with amino or hydroxyl groups on carbon atoms next to heteroaromatic nitrogens, ensuring precise SMILES strings to capture structure and functionality effectively.

# **Pinpointing amino positions**

OK, let me see. The amino group in 2-aminopyrimidine is located on carbon 2, with nitrogens at positions 1 and 3, ensuring precise SMILES notation for accurate molecular representation.

# Mapping substituent positions

I'm considering various patterns to place the amino group next to a pyrimidine ring nitrogen. Bracket notation and explicit substitutions are also being evaluated for accurate representation.

#### **Crafting unique structures**

I'm developing various bipyrimidine and pyrimidine-pyridine hybrids with specific substituent placements, aiming for precise SMILES strings to ensure clear molecular differentiation and correct nitrogen interaction.

#### **Crafting hybrid molecules**

I'm piecing together bipyrimidine and pyrimidine-pyridine hybrids, thinking through quinoxaline derivatives with precise substituent placements to ensure accurate SMILES notation and clear molecular differentiation.

#### Selecting substituent positions

I'm zeroing in on carbon 2 in quinoxaline to add an amino group, ensuring precise molecular representation in the final SMILES notation. This approach should streamline our tasks significantly.

#### Inserting the amino group

I'm piecing together the SMILES notation for 2-aminquinoxaline by carefully placing an NH2 group at carbon 2, ensuring precise molecular representation and enhancing clarity.

#### Generating new structures

I'm crafting isoquinoline derivatives with NH2 and OH groups at specific carbon positions, ensuring clear SMILES notation. Additionally, I'm exploring a bipodal structure with quinoxaline units.

#### Navigating molecular nuances

I'm thinking through bipodal quinoxaline structures with substituents on carbon near interacting nitrogens, including diamino, dihydroxy, and mixed forms. Exploring phenanthroline derivatives with NH2 groups for precise SMILES notation.

# **Ensuring accuracy**

I'm working through pyridazine derivatives, making sure substituents are positioned correctly. Checking SMILES strings for syntactic accuracy and suitable substituent placement near nitrogen atoms.

# **Ensuring precise bonding constraints**

I'm verifying that pyridazine derivatives have substituents on carbon atoms adjacent to nitrogen, avoiding quadruple bonds on nitrogen. I'll provide a list with each structure clearly labeled.

# **Enhancing structure clarity**

I'm mapping out clear SMILES notation for pyrimidine and pyridine derivatives, ensuring accurate substituent placement and avoiding ambiguity in chemical representations.

# Aligning substituent positions

I'm ensuring amino and hydroxy substituents are correctly positioned adjacent to nitrogen atoms in pyrimidine and pyridine derivatives, aiming for precise SMILES notation and clarity in representation.

Below is a list of 20 candidate molecules (provided as SMILES strings) in which the  $-NH_2$  or -OH substituents have been "moved" off the heteroatom and placed on an aromatic carbon adjacent to the heteroatom that is expected to interact with CO<sub>2</sub>. (In each design the "active" nitrogen in the ring remains unsubstituted so that its lone pair is available for binding, while an  $-NH_2$  or -OH group is attached on the neighboring carbon to provide a nearby hydrogen-bond donor and/or boost the electron density of the ring.)

**Important:** These SMILES strings represent "in silico" design hypotheses. Their chemical validity and performance must be confirmed by further (quantum-chemical or experimental) evaluation.

Below the candidate set is organized roughly by the core ring type.

# A. Pyridine-Based Systems

(Here the pyridine ring is defined with the ring nitrogen at position 1; an amino or hydroxy is placed at the adjacent C2.)

1. **Bipyridine Diamino** (two 2-aminopyridine units directly linked)

n1c(N)cccc1-c2c(N)cccc2

2. Bipyridine Dihydroxy (two 2-hydroxypyridine units)

nlc(O)ccccl-c2c(O)cccc2

3. Bipyridine Mixed (one ring 2-aminopyridine, the other 2-hydroxypyridine)

n1c(N)cccc1-c2c(O)cccc2

4. **Monocyclic 2-Aminopyridine** (as a reference; the free N at position 1 is available for binding)

n1c(N)cccc1

5. Monocyclic 2-Hydroxypyridine

n1c(O)cccc1

# **B.** Pyrimidine-Based Systems

(The pyrimidine ring here is drawn with nitrogens at positions 1 and 3. We assume that one of these N atoms serves as the  $CO_2$ -binding site; an  $-NH_2$  or -OH is attached at the C2 position – i.e. the carbon immediately adjacent to the N at position 1.)

# 6. Bipyrimidine Diamino

n1c(N)ncnc1-c2c(N)ncnc2

# 7. **Bipyrimidine Dihydroxy**

n1c(O)ncnc1-c2c(O)ncnc2

# 8. Bipyrimidine Mixed

n1c(N)ncnc1-c2c(O)ncnc2

# 9. Monocyclic 2-Aminopyrimidine

n1c(N)ncnc1

# 10. Monocyclic 2-Hydroxypyrimidine

n1c(O)ncnc1

# C. Mixed Heterocycle Systems

(Combining different rings to bring together two nitrogen-rich sites.)

11. **Pyrimidine–Pyridine Diamino** (left unit: 2-aminopyrimidine; right unit: 2-aminopyridine)

n1c(N)ncnc1-c2c(N)cccc2

12. **Pyrimidine–Pyridine Mixed (OH/NH<sub>2</sub>)** (left: 2-hydroxypyrimidine; right: 2-aminopyridine)

n1c(O)ncnc1-c2c(N)cccc2

13. **Pyrimidine–Pyridine Mixed (NH<sub>2</sub>/OH)** (left: 2-aminopyrimidine; right: 2-hydroxypyridine)

n1c(N)ncnc1-c2c(O)cccc2

# **D. Fused Bicyclic Heterocycles**

(These "condensed" systems often feature two rings sharing edges; substituents are placed on a ring carbon adjacent to a ring N.)

#### 14. 2-Aminoquinoxaline

(Quinoxaline core with the  $-NH_2$  at the C2 position; the interacting N (at position 1) remains free.)

c1nc2c(c1N)cccn2

#### 15. 2-Hydroxyquinoxaline

clnc2c(clO)cccn2

#### 16. Bipodal Quinoxaline Diamino

(Two quinoxaline units joined by a single bond; each carries an –NH<sub>2</sub> at the carbon adjacent to its "active" N.)

c1nc2c(c1N)cccn2-c3nc4c(c3N)cccn4

#### 17. Bipodal Quinoxaline Dihydroxy

c1nc2c(c1O)cccn2-c3nc4c(c3O)cccn4

#### 18. Bipodal Quinoxaline Mixed (NH<sub>2</sub>/OH)

c1nc2c(c1N)cccn2-c3nc4c(c3O)cccn4

# **E.** Other N-Rich Heterocycles

(Additional systems include isoquinoline, phenanthroline, or pyridazine derivatives.)

#### 19. 1-Aminoisoquinoline

(Isoquinoline with an  $-NH_2$  at the carbon adjacent to the ring N; here the substituent is placed on the benzene portion so that the ring N, typically at position 2 of isoquinoline, is free.)

Nc1ccc2nccc2c1

#### 20. 1-Hydroxypyridazine

(Pyridazine ring – here drawn as a six-membered ring with two adjacent N atoms; the – OH is placed at the C2 position, adjacent to the first ring N.)

n1c(O)cccn1

Each of these SMILES strings has been crafted so that no nitrogen in the heterocycle is "oversubstituted" (i.e. none bear four bonds) while still providing a neighboring  $-NH_2$  or -OH group as specified. You may now proceed with further quantum–chemical evaluation to assess their interaction energies with  $CO_2$ .

#### February 5, 2025 (ChatGPT-4o and reasoning)

#### User:

I am attaching a PDF with all the DFT-computed interaction energies from the molecules that you recommended based on our discussion. Based on these data, I would like from you to recommend 20 additional molecules (as SMILES strings) that have the potential to have interaction energy with CO2 that is lower than -10 kcal/mol. The main idea is to look for molecules with an electron rich site (e.g. N) that can induce a non-permanent dipole moment to the CO2, while additional H atoms from OH or NH2 group can further stabilize the CO2 molecule.

#### Response #11:

Based on your DFT-computed interaction energies, the best candidates have interaction energies below -10 kcal/mol, which means we need to explore structures that enhance electron-rich sites and hydrogen bonding potential.

Here are 20 additional molecules (in SMILES format) that could potentially show strong CO<sub>2</sub> affinity:

- 1. C1=NC2=C(N=CN=C2)N=C1OCCOCCO
- 2. C1=NC2=C(N=CN=C2)N=C1C3=NC=NC=C3O
- 3. C1=NC2=C(N=CN=C2)N=C1C3=NC=NC(=N3)O
- 4. C1=NC2=C(N=CN=C2)N=C1C3=NC=NC(=N3)N
- 5. C1=NC2=C(N=CN=C2)N=C1C3=NC=NS3
- 6. C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCO
- 7. C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCN
- 8. C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)NCCO
- 9. C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)NCCN
- 10. C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCO
- 11. C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCO
- 12. C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCOCCO
- 13. C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCOCCOCCO
- 14. C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCOCCOCCOCCO

These molecules **enhance electron density**, provide **multiple hydrogen bonding sites**, and **induce a dipole moment** in CO<sub>2</sub>, following the principles you outlined. Let me know if you'd like me to refine or filter based on specific criteria!

### S2. Prompts and Responses from Llama-3

# February 8, 2025 (Llama3)

# User:

*An active field of research is the search of molecules with strong interaction with carbon dioxide* (CO2). These CO2-philic units can be introduced in materials such as metal-organic frameworks or polymers, and through enhanced interactions with CO2, they can be used for the selective separation and capture of CO2. For example, Vogiatzis et al. estimated that the interaction energy between pyridine (SMILES: C1=CC=NC=C1) and CO2 is -4.46 kcal/mol, and between adenine (SMILES: C1=NC2=NC=NC(=C2N1)N) and CO2 is -5.90 kcal/mol. Similarly, Tian et al. found that the interaction energy between triazole-based molecules that can go up to -7.00 kcal/mol. More recently, Townsend et al. used machine learning and discover new organic molecules with enhanced CO2-philicity. For example, the molecule with SMILES string C1(=C(N=C([N]1[H])O[H])N([H])[H])N([H])[H] has a predicted interaction energy with CO2 of -7.41 kcal/mol. However, for enhance performance, it has been hypothesized that the interaction energy between these organic molecules and CO2 should be lower, at about -10 kcal/mol. I would like to ask you the following two questions: 1. Can you compile additional data from the literature with similar CO2-philic groups. You can either provide their SMILES string or their molecular formula. Listing the corresponding references would be much appreciated. 2. Based on the data that you have found and read, can you predict molecules with CO2 interaction that is lower than -7.5 kcal/mol and closer to -10 kcal/mol? Thanks

Note: Ollama was not able to recognize the links to the papers, so they were removed from the prompt

# Response #1:

I'd be happy to help!

1. Additional CO2-philic groups from literature:

- \* 4-Aminoazobenzene (SMILES: C12=CC=C(C=C1)N=Nc3ccc(cc3)NH2, reference: [1]) + Interaction energy with CO2: -5.15 kcal/mol
- \* Piperazine (SMILES: CN1C2CNCC2N1, reference: [2]) + Interaction energy with CO2: -4.42 kcal/mol
- \* Dicyclohexylurea (SMILES: CC1(CCC1)CCC1(NC(=O)NC1(CCC1)CCC1)N, reference: [3]) + Interaction energy with CO2: -6.13 kcal/mol
- \* 5-Aminouracil (SMILES: CN1CN=NC2=C(C=N1)N=C2, reference: [4]) + Interaction energy with CO2: -5.67 kcal/mol
- \* 3-Hydroxyindole (SMILES: O[C@H]1C2=CC(=C1)C=C2, reference: [5])
  - + Interaction energy with CO2: -6.14 kcal/mol

References:

[1] Vogiatzis et al., J. Phys. Chem. A 2018, 122(43), 8510-8523

[2] Tian et al., J. Phys. Chem. A 2020, 124(11), 2215-2226

[3] Li et al., ACS Sustainable Chem. Eng. 2020, 8(12), 4634-4644

[4] Wang et al., J. Mol. Struct.: THEOCHEM 2019, 1192, 126-135

[5] Zhang et al., RSC Adv. 2020, 10(31), 18333-18344

2. Predicting molecules with CO2 interaction energy lower than -7.5 kcal/mol and closer to -10 kcal/mol:

Based on the available data, I've identified some trends and features that could be useful for designing new CO2-philic molecules. These include:

\* Aromatic or heteroaromatic rings (e.g., pyridine, triazole) that provide a favorable environment for CO2 binding

\* Amide or urea groups (-NH-CO-NH-) that can form hydrogen bonds with CO2

\* Hydroxyl (-OH) or amino (-NH2) groups that can participate in hydrogen bonding with CO2 \* Cyclohexane or other cycloaliphatic rings that provide a hydrophobic environment for CO2

binding

Using these trends and features, I've generated some hypothetical molecules that could potentially have an interaction energy with CO2 lower than -7.5 kcal/mol and closer to -10 kcal/mol. Keep in mind that this is speculative, and actual performance may vary.

- \* Molecule 1: C12=C(C4=CC(=C3)N(C2=C[NH]N[H])O[C@H]1[C@H](C)(C3)=C4)C(F)(F) + Predicted interaction energy with CO2: -8.34 kcal/mol
- \* Molecule 2: CN1CN=NC2=C(C=N1)N=CC3=N[NH]O[C@H]1[C@H](C)(C3)=C4)C(F)(F) + Predicted interaction energy with CO2: -8.67 kcal/mol

\* Molecule 3:

CC1(CCCC[NH]O[C@H]1[C@H](C)(CC=C2=NC3=N[NH]N[H])O[C@H]1[C@H](C)(C3)=C 4)C(F)(F)

+ Predicted interaction energy with CO2: -9.14 kcal/mol

Please note that these predictions are based on simplified molecular models and may not accurately reflect the actual performance of these molecules in practice. Experimental validation is necessary to confirm their CO2-philicity.

References:

[1] Vogiatzis et al., J. Phys. Chem. A 2018, 122(43), 8510-8523
 [2] Tian et al., J. Phys. Chem. A 2020, 124(11), 2215-2226

Note: All SMILES provided by Llama were wrong, e.g.: SMILES Parse Error: unclosed ring for input: 'CC1(CCCC[NH]O[C@H]1[C@H](C)(CC=C2=NC3=N[NH]N[H])O[C@H]1[C@H](C)(C3)=C 4)C(F)(F)'

# **S3.** Prompts and Responses from Gemini 2.0

# February 8, 2025 (Gemini 2.0)

# User (Initial Prompt):

*An active field of research is the search of molecules with strong interaction with carbon dioxide* (CO2). These CO2-philic units can be introduced in materials such as metal-organic frameworks or polymers, and through enhanced interactions with CO2, they can be used for the selective separation and capture of CO2. For example, Vogiatzis et al. (ChemPhysChem 2009, 10, 374, https://doi.org/10.1002/cphc.200800583) estimated that the interaction energy between pyridine (SMILES: C1=CC=NC=C1) and CO2 is -4.46 kcal/mol, and between adenine (SMILES: C1=NC2=NC=NC(=C2N1)N) and CO2 is -5.90 kcal/mol. Similarly, Tian et al. (J. Phys. Chem. A 2015, 119, 3848, https://pubs.acs.org/doi/10.1021/acs.jpca.5b01892) found that the interaction energy between triazole-based molecules that can go up to -7.00 kcal/mol. More recently, Townsend et al. (Nat. Commun. 2020, 11, 3230, https://www.nature.com/articles/s41467-020-17035-5) used machine learning and discover new organic molecules with enhanced CO2philicity. For example, the molecule with SMILES string C1(=C(N=C([N]1[H])O[H])N([H])[H])N([H])[H] has a predicted interaction energy with CO2 of -7.41 kcal/mol. However, for enhance performance, it has been hypothesized that the interaction energy between these organic molecules and CO2 should be lower, at about -10 kcal/mol.

*I* would like to ask you the following two questions:

- 1. Can you compile additional data from the literature with similar CO2-philic groups. You can either provide their SMILES string or their molecular formula. Listing the corresponding references would be much appreciated.
- 2. Based on the data that you have found and read, can you predict molecules with CO2 interaction that is lower than -7.5 kcal/mol and closer to -10 kcal/mol?

Thanks :)

# Response #1:

You've asked a great question that delves into the cutting edge of CO2 capture research. Predicting and designing molecules with specific binding energies is a complex task, even with computational tools. I can't perform the level of *ab initio* calculations or advanced machine learning required for precise predictions. However, I can offer some strategies and point you towards relevant chemical motifs that are likely to exhibit strong CO2 interactions, and suggest how one might approach the design problem.

# 1. Compiling Additional Data on CO2-philic Groups:

It's important to note that finding *precise* interaction energies for a wide range of molecules with CO2 is challenging. Many papers report binding energies within specific contexts (e.g., embedded in a polymer, within a MOF) or use different computational methods, making direct

comparison difficult. However, we can look at functional groups known to interact favorably with CO2:

- Amines (Primary, Secondary, Tertiary, and especially Sterically Hindered Amines): These are the workhorses of CO2 capture. Their interaction with CO2 leads to the formation of carbamates, bicarbonates, and carbamic acids. The strength of the interaction depends heavily on the amine's structure and environment. Many variations exist, and research continues to optimize them. It's difficult to give a single SMILES string as the possibilities are vast. Look for papers on "amine-functionalized CO2 capture materials" for numerous examples.
- **Guanidines:** These functionalities exhibit strong basicity and can form strong interactions with CO2. They are less common than amines but hold promise. A simple guanidine could be represented as C (N) (N) =N.
- Amides: While not as strong as amines, amides can still contribute to CO2 binding, especially when multiple amides are present in a molecule. C (N) =0 is the basic amide functionality.
- Ureas: Similar to amides, ureas can interact with CO2. NC(N) = 0 is the urea functionality.
- Heterocyclic Nitrogen-containing rings: Pyridine, imidazole, triazole, and similar rings (as you mentioned) can contribute to CO2 binding through the nitrogen atom. Their interaction energies vary based on the specific ring and substituents.
- **Boron-containing groups:** Boron atoms can also interact with CO2. For example, boronic acids and their derivatives.

# References (General Areas - Specific values are hard to pinpoint without context):

- *Chemical Reviews* often has review articles on CO2 capture materials, which will contain many examples.
- Journal of the American Chemical Society, Angewandte Chemie, Chemical Communications, Energy & Environmental Science are good places to search for primary research articles.

# 2. Predicting Molecules with Enhanced CO2 Interaction:

Reaching -10 kcal/mol is a significant challenge. Here's a strategy combining chemical intuition and computational considerations:

- **Combine Multiple CO2-philic Groups:** The key is likely to combine several of the functional groups mentioned above within a single molecule. For instance, a molecule with multiple sterically hindered amines *and* some strategically placed heterocyclic nitrogen rings could exhibit a synergistic effect.
- Steric Hindrance: In the case of amines, steric hindrance around the nitrogen atom can actually *increase* the CO2 binding strength by preventing the formation of less stable carbamate species and favoring bicarbonate formation.
- **Computational Screening:** The best approach would be to use computational chemistry (DFT calculations, for example) to screen a large library of candidate molecules. This is

what the Townsend et al. paper you mentioned did, using machine learning to guide the search. You would need access to these computational resources and expertise.

- Focus on Rigidity: Rigid molecules (less flexible) can sometimes exhibit stronger binding because they don't lose as much entropy upon binding.
- **Consider Cooperative Binding:** Designing molecules where multiple binding sites work together cooperatively could enhance the overall binding energy.

# **Example (Conceptual - Requires Computational Validation):**

Imagine a molecule with a rigid scaffold (perhaps a polycyclic structure) to which you attach multiple sterically hindered amines (like diisopropylamine or similar) and some triazole or imidazole groups. The idea is to maximize the number of CO2 interaction sites in a constrained environment. However, this is just a starting point. Computational chemistry is essential to refine the design and predict the actual binding energy.

**Important Note:** I cannot provide specific SMILES strings with guaranteed -10 kcal/mol binding energy. Accurate prediction requires advanced computational methods. My suggestions aim to guide you towards promising chemical motifs and design strategies. Collaboration with computational chemists is crucial for this type of research.

# <u>User:</u> Could you type such molecules in SMILES format?

# Response #2:

I understand you'd like SMILES strings for potential CO2-philic molecules. However, I must emphasize that I cannot *design* molecules with specific binding energies. Predicting binding energies requires complex quantum chemical calculations, not just stringing together SMILES. Providing SMILES strings without that computational backing would be misleading.

What I *can* do is give you SMILES examples of *parts* of molecules that are known to be CO2philic. You can then use these as building blocks in your own design process, keeping in mind that the overall molecule's structure and environment will drastically affect the actual CO2 interaction. *These examples do not guarantee any specific binding energy*.

# Examples of CO2-philic fragments (not complete molecules):

- Sterically hindered amine: CC(C)(C)NC(C)(C)C(Diisopropylamine a *part* of a larger molecule)
- **Triazole:** C1=CN=NC=C1
- Imidazole: C1=CN=C1
- Guanidino group: C(N) (N) = N (Again, this needs to be attached to something)
- Urea: NC (N) =0

# Why I can't give complete "designed" molecules:

- 1. **Computational Complexity:** Accurately predicting CO2 binding energies requires Density Functional Theory (DFT) calculations or other sophisticated quantum chemical methods. I'm a language model; I can't perform those calculations.
- 2. **Combinatorial Explosion:** Even for a relatively small molecule, the number of possible arrangements of CO2-philic groups is enormous. It would be impossible to simply guess and hope for a specific binding energy. Systematic exploration requires computational screening.
- 3. **Context Matters:** The CO2 binding energy isn't just a property of the isolated molecule. It's also affected by the molecule's environment (e.g., if it's in a polymer matrix, a MOF, or in solution). I can't account for those environmental effects in a SMILES string.

# What you should do:

- 1. Learn Computational Chemistry: If you want to design molecules with specific properties, you'll need to learn computational chemistry methods (or collaborate with someone who does).
- 2. Use Molecular Modeling Software: Software packages like Gaussian, ORCA, or Q-Chem are used for these calculations.

- 3. **Start with Simple Systems:** Begin by studying simpler molecules to understand the principles before tackling complex designs.
- 4. **Consult the Literature:** Look for papers that report CO2 binding energies for specific molecules. This will give you a better idea of what's possible and how different functional groups contribute.

I hope this explanation is helpful. I want to be transparent about my limitations and guide you towards the right resources for this type of research.

# User:

Could you combine these ideas that you just presented and provide molecules as SMILES strings with potential string CO2 interactions (about -10 kcal/mol)?

# Response #3:

While I can't guarantee these molecules will have exactly -10 kcal/mol interaction energy (as that requires computational chemistry and validation), here are some SMILES strings combining the ideas discussed, with the caveat that these are conceptual and would need further study:

# 1. Multiple Amine and Steric Hindrance:

- C1 (CN (C1) C (C (C) C) C (C (C) C) N (C (C) C) C (C) C) N (C (C) C) C (C) C
  - This has a central piperazine core with multiple bulky isopropyl groups attached to the nitrogens, creating steric hindrance around the amine binding sites.

# 2. Guanidine with a Hindered Amine:

- C(NC(N)=NC1CC2CC1CN2)(NC(C)C)C(C)C
  - A guanidine group is attached to a cyclobutane ring, and a sterically hindered isopropylamine is also attached to the central carbon.

# 3. Macrocycle with Amine and Urea:

- C1 (=0) NC (NC1) (C2CN (C2) C3CC4CC3CN4) N
  - A urea group is incorporated into a macrocycle along with a piperazine-like unit. The macrocycle could provide a pre-organized binding site for CO2.

# 4. Triazole with Multiple Amines:

- C1=NC=NC(=C1N(C(C)C)C(C)C)N(C(C)C)C(C)C
  - A triazole ring is substituted with two sterically hindered amine groups.

# 5. Cage-like Structure with Guanidines:

- C1CN2C1CN(C2)C3CN4C3CN(C4)C(=NC(N)=N)(NC(C)C)C(C)C
  - This is a simplified representation of a cage-like structure where guanidine groups are attached to the cage framework. (Synthesizing such a molecule would be a significant challenge).

# **Important Considerations:**

• **Synthesis:** These are just SMILES strings. Synthesizing these molecules could be very complex and might require multiple steps.

- **Computational Validation:** It's essential to use computational chemistry methods (DFT, etc.) to estimate the CO2 binding energies of these molecules. I cannot guarantee they will be near -10 kcal/mol.
- **Experimental Validation:** Even if computational studies are promising, experimental validation is crucial. The actual CO2 binding strength in a real-world setting might differ from theoretical predictions.
- **Optimization:** These are just starting points. Further optimization of the structures (e.g., by adding different substituents) might be needed to achieve the desired CO2 binding energy.

Remember, designing molecules for CO2 capture is an iterative process that involves computational design, synthesis, characterization, and optimization. I hope these ideas inspire you to explore further!

# User:

Thank you! I got an error about the last SMILES that you sent me: Could not parse SMILES C1CN2C1CN(C2)C3CN4C3CN(C4)C(=NC(N)=N)(NC(C)C)C(C)CExplicit valence for atom # 14 C, 5, is greater than permitted Can we fix this?

# Response #4:

You are absolutely right! My apologies. I made a mistake in the valence of one of the carbons in the cage structure. A carbon atom can only have a maximum of four bonds. I was trying to represent a more complex cage structure conceptually, and in doing so, I violated this rule.

Here's a corrected and slightly simplified version of the cage-like structure with guanidines, and a few other variations to explore. It's still a simplified representation, and synthesizing such a molecule would be a significant synthetic challenge, but at least the SMILES strings are now valid:

# Corrected and Simplified Cage-like Structures:

- 1. C1CN2C1CN (C2) C3CN4C3CN (C4) C (=NC (N) =N) N (This version has one less attached group, and it is synthetically simpler. It still contains a cage-like structure with guanidine group.)
- 2. C1CN2C1CN (C2) C (=NC (N) =N) N (Even simpler cage-like structure with guanidine group)

# **Other related structures:**

- 3. C1CC (NC1) (C (=NC (N) =N) N) C2CC2 (A bicyclic structure with a guanidine)
- 4. C1 (CN2C1CN2) C (=NC (N) =N) N (A simpler bicyclic structure with a guanidine)

# **Explanation of the Error and Correction:**

The original SMILES string I provided had a carbon atom (atom #14 in the error message) that was bonded to five other atoms. This is impossible for carbon. The corrected versions maintain the cage-like or bicyclic structure but reduce the number of substituents on some of the carbon atoms, ensuring that no carbon has more than four bonds.

Important Note: Even with the corrected SMILES strings, the same caveats apply:

- Synthesis: Synthesizing these molecules would be a significant synthetic challenge.
- **Computational Validation:** It's essential to use computational chemistry methods to estimate the CO2 binding energies.
- **Experimental Validation:** Experimental validation is crucial.
- **Optimization:** These are just starting points.

I recommend using a molecule drawing program (like ChemDraw, Marvin Sketch, or RDKit) to visualize these structures. It might help you to understand the connectivity better. Also, remember that even if the SMILES string is valid, it doesn't guarantee that the molecule is stable or even synthesizable. Computational and experimental work would be needed to assess that.

#### **S4.** Computational Details

The computational workflow used in this study is shown graphically in Figure 1 of the main manuscript. Here, we present the details of each of the individual steps.

### 1. SMILES Conversion to Cartesian Coordinates

All SMILES generated by GPT-40 and Gemini 2.0 were converted into Cartesian coordinates with the RDKit,<sup>1</sup> an open-source toolkit for cheminformatics, and with the Chem.MolFromSmiles module. Synthesizability scores (synthetic accessibility score,<sup>2</sup> natural product-likeness score<sup>3</sup>), as implemented on RDKit, were used for selected molecules.

# 2. Automated Conformer Analysis

Due to the structural complexity of most of AI-generated molecular structures, and to ensure a fully automated and unbiased geometry optimization of both the molecular unit and the unit-CO<sub>2</sub> supersystem, we have utilized the Conformer–rotamer sampling tool (CREST).<sup>4</sup> We executed CREST with the GFN2-xTB<sup>5</sup> semi-empirical tight-binding method and a Fermi temperature of 300K. The 10 most stable conformers were selected for further optimization with density functional theory (DFT) calculations. In a few cases, less than 10 conformers within an energy threshold of 6 kcal/mol were obtained from CREST and thus, less than 10 DFT geometry optimizations were performed. For the conformer sampling of the CO<sub>2</sub>-organic molecule supersystem, we use the NCI mode of CREST that enables to find aggregates of noncovalent complexes.

# 3. Density Functional Theory Calculations

All DFT calculations were performed with the TURBOMOLE 7.2.1 quantum chemical program package. The PBE0 density functional<sup>6</sup> was used in all calculations together with Grimme's D4 dispersion correction,<sup>7</sup> the def2-TZVPP basis set,<sup>8</sup> the resolution of identity (RI) approximation,<sup>9</sup> and the m4 integration grid. All geometry optimizations were performed with strict convergence criteria (convergence total energy and maximum norm of cartesian gradient were set to 10<sup>-8</sup> and 10<sup>-4</sup>, respectively). The PBE0 functional was selected based on previous benchmark studies that has shown that this functional together with triplezeta quality basis sets provide high accuracy close to CCSD(T) at the complete basis set limit.<sup>10, 11</sup> The energetically most stable conformer for each molecular species (organic molecule or molecule-CO<sub>2</sub> supersystem) was selected for further analysis. Vibrational frequency calculations were performed on the most stable conformations based on the DFT geometry optimizations to ensure they are minima on the potential energy surface. No corrections for the basis set superposition error (BSSE) were considered in this study since a significantly large basis set was applied in all DFT calculations (def2-TZVPP). In a previous study,<sup>11</sup> we found that the that PBE0-D3(BJ) with def2-TZVPP has a deviation of only 0.13 kcal/mol from CCSD(F12)(T)/CBS.

# 4. Interaction Energies

Interaction energies ( $\Delta E_{Int}$ ) between two molecular unit and CO<sub>2</sub> were calculated as:

$$\Delta E_{\rm Int} = E_{\rm mol-CO2} - E_{\rm mol} - E_{\rm CO2}$$

where  $E_{mol-CO2}$  represents the energy of the weakly interacting supersystem, and  $E_{mol}$  and  $E_{CO2}$  represent the energies of the isolated, noninteracting molecules. Solvent effects were not considered in this study since the focus of the CO<sub>2</sub>-philicity is on solvent-free technologies, such as physisorption or diffusion-based separations. Condense-phase effects were excluded as we rely on the fact that the CO<sub>2</sub>-philicity is an indicator of the physisorption or diffusion.

# **S5. List of Molecular Units and Computed Interaction Energies**







022a



,----















031a



031b

031c







031d

031e





































































чH

















× 100 Ľγ

the state



Ser P



















GI 1-	40	
SMILES	Note	$\Delta \overline{E_{\text{Int}}}$
01 N(CCN(CCN)C=O)C=O		-4.35
02 NCCOCCOCCOCCN		-5.25
03 C1(C(=O)NC2=CC=C(C=C2)C(=O)NCC(=O)NCC(=O)NCC(=O) (=O)NCC(=O)NC1)	D)NCC(=O)NCC	-5.28
04 N1C=CN=C1CCOCCOCCOCCN1C=CN=C1		-5.59
<b>05</b> C(CN(CC(=0)O)CC(=O)O)(CN(CC(=O)O)CC(=O)O)N		-8.15
06 C10C2(C01)OC(C0)C(O2)N		-4.95
07 N(CCOCCOCC(=0)NC(=0)OCCOCCN)C(=0)NC(=0)O		-6.24
08 N1C=CN=C1CC(CN(CCOCCOCCN)CCN1C=CN=C1)C(=O)O		-4.17
09 N1(C(CCOCCOCCNCCOCCNCCOCCO)C1)	Corrected by GPT-40	-5.31
10 N1C=CN=C1CCOCCOCCOCC(=O)N(CCN1C=CN=C1)C	(=O)O	-5.53
11 N(CCOCCN(CC(=0)NC(=0)OCCOCCN)C(=0)O)C(=0)NC(=0)	0)0	-6.00
12 C1(=CC(=O)N(CCN1C=CN=C1)C(=O)N(CCOCCOCCN)C(=O)	)O)N1C=CN=C1 Erroneous SMILES - excluded	-
13 N(CCN(CCN(CC(=0)NC(=0)N1C=CN=C1)C(=0)O)C(=0)O)C	C(=O)O	-6.06
14 C1=C(N=CN=C1O)C2=NC=NC(=C2O)C3=NC=NC=C3O		-5.27
15 C1=NC(=NC=N1)C2=NC=CN=C2OCCO		-4.62
16 C1=CC2=NN=CN2C=C10CCOCCO		-6.54
17 C1=NC(=NO1)C2=NC=NC=N2C3=NC=NO3	Corrected manually	-4.77
18 C1=NC2=CC=NC=C2N=C1C3=NN=NN3C(=O)O		-4.09
19 C1=NC=CN=C1C2=CN=NN2C3=NC=CS3		-5.08
20 C1=NC(=NC(=N1)C2=NC3=CC=CC=C3N2)OCCO		-5.82
21 C1=NC2=CC=NC=C2C3=NN=CN=C31		-4.11
22 C1=NC2=C(N=CN2C=N1)OCCOCC3=NC4=C(N=CN4C=N3)C	)	-6.44
<b>22-a</b> C1=NC2=C(N=CN2C=N1)OCCOCC3=NC4=C(N=C(N)N4C=N	(3)O Added by user	-7.80
<b>22-b</b> C1=NC2=C(N=CN2C=N1)OCCOCC3=NC4=C(N=C(O)N4C=N	(3)O Added by user	-7.91
23 C1=NC2=C(N=CN2C=N1)C3=NN=NN3C(=O)O		-5.05
24 C1=NC2=C(N=CN2C=N1)C3=NC4=CC=NC=C4N3		-5.73
25 C1=NC2=C(N=CN2C=N1)N3C=NC4=CC=CC=C43OCCO	Corrected manually	-5.71
26 C1=NC2=C(N=CN2C=N1)OCCOCCOCCOCCO		-5.86
27 C1=NC2=C(N=CN2C=N1)C3=NC4=C(N=CN4C=N3)C5=NC6= OCC(=Q)Q	=C(N=CN6C=N5)	-7.04
<b>28</b> C1=NC2=C(N=CN2C=N1)C3=NC=NS3		-5.01
29 C1=NC2=C(N=CN2C=N1)OCCOCCOCCOCCOC		-6.20
<b>30</b> n1c(N)nc(c1)-c2nc(N)nc(n2)		-6.24
<b>31</b> n1c(O)nc(c1)-c2nc(O)nc(n2)		-7.21
<b>31-a</b> n1c(O)nc(c1)-c2nc(O)cc(n2)	Added by user	-7.39
<b>31-b</b> n1c(O)nc(c1)-c2nc(O)nc(n2)N	Added by user	-7.27
<b>31-c</b> n1c(O)nc(c1)-c2nc(O)cc(n2)N	Added by user	-7.41
<b>31-d</b> n1c(O)nc(c1)-c2nc(N)cc(n2)N	Added by user	-7.42
<b>31-e</b> n1c(O)nc(c1N)-c2nc(O)nc(n2)	Added by user	-7.23
32 n1c(N)nc(c1)-c2nc(O)nc(n2)		-6.30
33 c1c(nccn1N)-c2c(nccn2N)	Corrected manually	-5.29
<b>34</b> c1c(nccn1O)-c2c(nccn2O)	Corrected manually	-5.11
<b>35</b> Nc1nc(N)nc(N)c1-c2nc(N)nc(N)c2		-6.14

36	Nc1nc(N)nc(N)c1-c2nc(N)nc(O)c2		-6.12
37	n1c(N)nc(c1)-c2ncccn2	Corrected by GPT-40	-5.90
38	n1c(O)nc(c1)-c2ncccn2	Corrected by GPT-40	-6.80
39	n1c(N)nc(c1)-c2ncccn2O	Corrected manually	-6.28
40	n1c(NC)nc(c1)-c2nc(NC)nc(n2)	Corrected by GPT-4o	-6.78
41	n1c(N)nc(c1)-c2nc(N)nc(c2)N	Corrected by GPT-40	-6.07
42	c1ncnc2ncn(N)c2n1	Corrected by GPT-40	-4.82
43	c1ncnc2ncn(O)c2n1	Corrected by GPT-40	-4.74
44	c1ncnc2n(cn1)nc(N)c2		-5.34
45	Nc1ncncn1-c2ncncn2	Corrected manually	-5.32
46	Oc1nenen1-c2nenen2	Corrected manually	-5.25
47	n1c(N)nc(c1)-c2ncncn2	Corrected by GPT-40	-5.63
48	n1c(O)nc(c1)-c2ncncn2N	Corrected manually	-5.51
49	n1c(N)nc2c1ncnc2O		-4.75
50	n1c(N)cccc1-c2c(N)cccc2		-5.32
51	n1c(O)cccc1-c2c(O)cccc2		-5.25
52	n1c(N)cccc1-c2c(O)cccc2		-4.90
53	n1c(N)cccc1		-5.51
54	nlc(O)ccccl		-5.44
55	n1c(N)ncnc1-c2c(N)ncnc2		-5.20
56	n1c(O)ncnc1-c2c(O)ncnc2		-3.98
57	n1c(N)ncnc1-c2c(O)ncnc2		-5.92
58	n1c(N)ncnc1		-5.18
59	n1c(O)ncnc1		-5.22
60	n1c(N)ncnc1-c2c(N)cccc2		-5.31
61	n1c(O)ncnc1-c2c(N)cccc2		-5.44
62	n1c(N)ncnc1-c2c(O)cccc2		-5.17
63	clnc2c(clN)cccn2		-5.77
64	c1nc2c(c1O)cccn2		-5.99
65	c1nc2c(c1N)cccn2-c3nc4c(c3N)cccn4		-5.77
66	c1nc2c(c1O)cccn2-c3nc4c(c3O)cccn4		-5.90
67	c1nc2c(c1N)cccn2-c3nc4c(c3O)cccn4		-5.90
68	Nc1ccc2nccc2c1		-4.88
69	nlc(O)cccnl		-5.31
70	C1=NC2=C(N=CN=C2)N=C1OCCOCCO		-3.82
71	C1=NC2=C(N=CN=C2)N=C1C3=NC=NC=C3O		-4.16
72	C1=NC2=C(N=CN=C2)N=C1C3=NC=NC(=N3)O		-6.31
73	C1=NC2=C(N=CN=C2)N=C1C3=NC=NC(=N3)N		-6.32
74	C1=NC2=C(N=CN=C2)N=C1C3=NC=NS3		-4.41
75	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCO		-4.00
76	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCN		-4.90
77	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)NCCO		-6.31
78	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)NCCN		-6.33
79	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCO		-8.04
80	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCO		-5.89

81	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCOCCO		-2.89
82	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCOCCOCCO	Ethylene oxide oligomer - excluded	-
83	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCOCCOCCOCCO	Ethylene oxide oligomer - excluded	-
84	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCOCCOCCOCCO CCO	Ethylene oxide oligomer - excluded	-
85	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCOCCOCCOCCO CCOCCO	Ethylene oxide oligomer - excluded	-
86	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCOCCOCCOCCO CCOCCOCCO	Ethylene oxide oligomer - excluded	-
87	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCOCCOCCOCCO CCOCCOCCOCCO	Ethylene oxide oligomer - excluded	-
88	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCOCCOCCOCCO CCOCCOCCOCCOCCO	Ethylene oxide oligomer - excluded	-
89	C1=NC2=C(N=CN=C2)N=C1C3=NC(=NC=N3)OCCOCCOCCOCCOCCOCCOCCOCCOCCOCCOCCOCCOCC	Ethylene oxide oligomer - excluded	-

	Gemini 2.0		
	SMILES	Note	$\Delta E_{\rm Int}$
90	C1(CN(C1)C(C(C)C)C(C(C)C)N(C(C)C)C(C)C)N(C(C)C)C(C)C		-1.23
91	C(NC(N)=NC1CC2CC1CN2)(NC(C)C)C(C)C		-3.96
92	C1(=0)NC(NC1)(C2CN(C2)C3CC4CC3CN4)N		-3.38
93	C1=NC=NC(=C1N(C(C)C)C(C)C)N(C(C)C)C(C)C		-3.24
94	C1CN2C1CN(C2)C3CN4C3CN(C4)C(=NC(N)=N)N		-7.24
95	C1CN2C1CN(C2)C(=NC(N)=N)N		-4.73
96	C1CC(NC1)(C(=NC(N)=N)N)C2CC2		-4.57
97	C1(CN2C1CN2)C(=NC(N)=N)N		-5.09



**Figure S1.** Distribution of computed interaction energies for the initial 88-molecule dataset and the 88+7 molecule dataset.

A small number of SMILES had errors that were either fixed by new prompts to the LLM, or manually by the user, except structure 12, that was excluded from further analysis. Structures 79 - 89 were including N-containing units together with linear chains with carbonyl groups of increasing size (from three up to 13 units). In most of these cases, the linear chains are blocking the CO<sub>2</sub>-philic sites due to favorable intramolecular interactions. For that reason, only structures with short chains were kept in our study and structures 82 - 89 were excluded. The final number of structures considered in the first stage of the study was 88. The molecular list was augmented by 7 additional structures that were added by the user and thus, the full dataset was extended to 95 molecules. Analysis of the initial 88 structures revealed a wide range of interaction energies between CO<sub>2</sub> and the molecular units. The average interaction energy of these 88 structures is -5.39 kcal/mol, while the majority of the suggested molecules were between -4.5 and -7 kcal/mol (78% out of 88 structures). The molecule with the least attractive interaction energy included two diisopropylamino groups (structure 90), where the terminal methyl groups are unable to induce a non-permanent dipole to CO<sub>2</sub>, limiting the capacity to attract the linear molecule and thus, leading to an interaction energy of  $\Delta E_{\text{Int}} = -1.23$  kcal/mol. From the remaining cases, five molecules (05, 27, 31, 79, 99, Figure 2) showed stronger CO<sub>2</sub> affinity with interaction energies exceeding the -7 kcal/mol.

	R <sub>C</sub>	<b>R</b> <sub>01</sub>	<b>R</b> <sub>01</sub>	$\Delta E_{\text{Int(CO2)}}$	$\Delta E_{\text{Int(N2)}}$	$\Delta E_{\text{Int(CH4)}}$	SA	NP
05	2.763	2.050	2.124	-8.15	-3.27	-3.09	2.688	-0.217
	2.915				(2.5)	(2.6)		
					` '	` <i>´</i>		
22	2.878	2.091	3.159	-6.44	-1.43	-1.58	3.458	-1.029
	2.977				(4.5)	(4.1)		
22-а	2.767	2.375	2.567	-7.80	-2.04	-2.70	3.534	-1.029
	3.034				(3.8)	(2.9)		
22-ь	2.873	2.003	3.221	-7.91	-1.51	-1.69	3.555	-1.069
	2.965				(5.2)	(4.7)		
27	2.867	2.587	-	-7.04	-3.51	-3.74	3.609	-0.712
					(2.0)	(1.9)		
31	2.893	1.976	3.020	-7.21	-2.90	-2.68	3.738	-0.478
	3.868				(2.5)	(2.7)		
<b>31-</b> a	2.895	1.987	2.849	-7.39	-2.97	-2.70	3.480	-0.917
	3.828				(2.5)	(2.7)		
31-ь	2.896	1.972	2.995	-7.27	-2.92	-2.69	3.641	-0.676
	3.855				(2.5)	(2.7)		
1	2.005	1.070	2.0.00	7.41	2.07	0.70	2 5 5 7	0.651
31-c	2.895	1.979	2.868	-7.41	-2.96	-2.70	3.557	-0.651
	3.829				(2.5)	(2.7)		
21 -1	2 000	1.077	2 000	7.40	2.00	2.75	2 274	0 (20
31-0	2.890 2.825	1.9//	2.909	-/.42	-2.89	-2.13	3.3/4	-0.020
	3.833				(2.0)	(2.7)		
31 0	2 871	1 076	3 105	7 72	2 87	2.68	3 707	0.512
31-6	2.0/1	1.7/0	5.105	-1.23	-2.07	-2.00	5.707	-0.312
	3.703				(2.3)	(2.7)		
79	2 871	2 136		-8.04	-2.66	-3 30	2 963	-0.880
17	2.071	2.130	-	-0.04	(3.0)	(2.4)	2.705	-0.000
	2.702				(3.0)	(2.7)		
94	2.924	2.688	2,740	-7.24	-2.39	-3.11	4.809	-0.241
<i>·</i> ·	3.087	2.951	,	,	(2.8)	(2.3)		U
	2.007	, 0 1			()	()		

# S6. Additional Analysis of Best Performers

Key:	

$R_{\rm C}$ :	Closest atom distances (in Å) from carbon of $CO_2$
$R_{O1}$ :	Closest atom distances (in Å) from $oxygen(1)$ of $CO_2$
$R_{O2}$ :	Closest atom distances (in Å) from $oxygen(2)$ of $CO_2$
$\Delta E_{\text{Int(CO2)}}$ :	Interaction energy (in kcal/mol) of molecular unit with CO <sub>2</sub>
$\Delta E_{\text{Int(N2)}}$ :	Interaction energy (in kcal/mol) of molecular unit with N <sub>2</sub>
	The ratio between $\Delta E_{\text{Int}(\text{CO2})}$ and $\Delta E_{\text{Int}(\text{N2})}$ is shown in parenthesis
$\Delta E_{\text{Int(CH4)}}$ :	Interaction energy (in kcal/mol) of molecular unit with CH <sub>4</sub>
	The ratio between $\Delta E_{\text{Int(CO2)}}$ and $\Delta E_{\text{Int(CH4)}}$ is shown in parenthesis
SA:	Synthetic Accessibility Score, <sup>2</sup> ranges from 1 to 10, 1 refers to easier synthesis
NP:	Natural Product-likeness Score, <sup>3</sup> ranges from -5 to 5, -5 refers to easier synthesis

#### **S7. Reproducability Experiments**

A known limitation of current LLMs is their reproducibility. Many components affect the final output, including prompt selection, version, and date. To explore this behavior, we have tested if ChatGPT-4o has the ability to reproduce the previous results. We have asked the same 5 initial prompts used in this study (see section S1) and we used ChatGPT-4o without logging in the user account. This process was repeated 4 times (termed as experiments in the discussion). From the first two prompts, the majority of the initial molecular units were different in the four experiments, except for triazine variants that were present in each attempt. As more questions were asked during the partial training of the LLM, we identified that all models start to converge to the molecular units that were discovered by the initial study. For example, molecular units that are composed by two bonded heteroaromatic groups (structure **31-d** from Figure 1 of manuscript) were present multiple times in each of the four experiments:

- Experiment 1: structures 19, 22, 25, 33, 34, 35
- Experiment 2: structures 23, 24
- Experiment 3: structures 10, 22, 25, 27, 28
- Experiment 4; structures 22, 23, 24, 25, 26

Similarly, molecular units with chains containing carboxylic units (structures **27** and **79** from Figure 1 of manuscript) were also present in each of the four experiments:

- Experiment 1: structures 32
- Experiment 2: structure 6
- Experiment 3: structures 32, 34, 40
- Experiment 4; structures 17, 18

Many structures from these four experiments also resembled the  $\pi$ - $\pi$  stacked geometry of **22-b**. Additional model training based on user's experience can further tune these molecules from the four experiments, similar to the work presented in the manuscript. Alternatively, statistical analysis by output aggregation from a large number of experiments can limit the diversity of the LLM's responses and support reproducibility in a systematic manner.

Overall, the following results were collected (incorrect SMILES were excluded):

# May 23, 2025 (ChatGPT-40)

#### Experiment 1:

001 N=C(N)N1C=CN=N1 002 N1=NC(N)=NN1N 003 Nc1nc(N)c(N)c(n1)O 004 NC(=N)NC(=N)N 005 Nc1ncnc(N)n1 006 O=c1[nH]cnc1O 007 NC(=N)NCO 008 NCCN1C=CN=N1 009 NC(=N)NC1=CN=NN1CO 010 NC(=N)NCC1=CN=NN1N 011 NC(=N)NNC=NO

```
012 N=C(N)NC1=NC(N)=NC(N)=N1
013 C1=NN=CN1COC(=N)N
014 NC(=N)NCC1=CN=NN1COCO
015 NC(=N)NCC1=NN=CN1CC2=CN=NN2
016 C12=NN=CN1CN=C(N)N2
017 NCC1=NN=CO1NCCN
018 NC(=N)NCC1=NC2=CC=CC=C2N1CO
019 OC1=NN=CN1C2=NN=CO2
020 O=C1ON=NC2=CN=C12
021 C1=CC2=NON=C2C(=N1)N=[N+]=[N-]
022 C1=CN=NN1C2=NO[N]O2
023 C1=CC2=NC=NO2C3=NOCO3
024 O1C=NN=C1C2=CN=CN2C3=NN=CO3
025 O=C1C=CN=CN1C2=NC=NO2
026 C1=CN2C=NC=CC2=N1C3=NO[N]O3
027 C1=CN=CC2=C1C=NC=N2C3=NOCO3
028 C1=CC2=NC=CN=C2C(=N1)C3=NOCO3
029 NC(=N)NCC1=NC2=C(N1)N=CN=C2O
030 OCNCC1=NC2=C(N1)N=CN=C2CO
031 OC(=O)C1=NC2=C(N1)N=CN=C2
032 OCCOCC1=NC2=C(N1)N=CN=C2
033 C1=NC2=C(N1)N=CN=C2C3=NOCO3
034 C1=NC2=C(N1)N=CNC2C3=NO[N]O3
035 C1=NC2=C(N1)N=CN=C2C3=CN=NN3
036 OC1=NC2=C(N1)N=C(O)N=C2O
037 OP(=O)(O)CC1=NC2=C(N1)N=CN=C2
038 O=S(=O)(O)CC1=NC2=C(N1)N=CN=C2
```











н н

H<sup>\*</sup>





















N H















032





024

028

030





033



034

029



031







037



# **Experiment 2:**

036

001 NC1=NC(N)=NC(N)=N1 002 N=C(N)NO 003 NC1=NC=CN1 004 NC(=N)NC(=N)N 005 NC1=NN=CN1O 006 C(C(=O)O)N=C(N)N 007 NC1=NN=CN1CNC2=NN=CN2N 008 NC1=NC(N)=NN1 009 N=C(N)C(=N)N 010 C(N1C=NN=C1)(N2C=NN=C2)N3C=NN=C3 011 C1(=NC(N)=N)NC(=N)N1 012 N=C(N)NC1=NN=CN1N 013 C1=NC2=NC(N)=NN2C(N)=N1 014 C1(NC(=N)N)=NN=CN1NC(=N)N 015 C1=NC(=O)N=N1C2=NN=CN2N 016 N1C=NN=C1CCN2C=NN=C2N 017 N=C(N)NCC(N=C(N)N)N 018 N1C=NC(N)=NC2=NN=CN12 019 C1=NN=CN1C2=NN=NC2 020 C#CC1=NN=CN1C2=NN=CN2 021 C1=NN=CN2C=NC=NC12















































#### **Experiment 3:**

001 C1=NN=C(N1)N=C(N)N 002 C1=NN=CN1C(=NH)NH2 003 NC1=NC(N)=NC(N)=N1 004 NC1=NC(NC2=NN=CN2)=NC(N)=N1 005 OC1=NN=CN1 006 N[C@@H](CO)C1=NN=CN1 007 N=C(N)NCCN=C(N)N 008 NC(N)=NC1=CC=CC=C1 009 C1=CC(=CC=C1C(=O)NCC(=NH)NH2)N2C=NC=N2 010 C1=CC=C(C=C1)C(NC(=O)N)NC2=NN=CN2 011 N=C(N)NCC(=O)NC(N)=NC(=NH)N 012 NCC(=O)NC1=NN=C(N1)C(=O)NCC(=NH)N 013 C1=CC=C(C=C1)C(=O)NC2=NN=C(N2)N=C(N)N 014 C1=CC=CC(=C1)NC(=O)NCC2=NN=C(N2)NCC(=O)N 015 N=C(N)NC(=O)NCC(=O)NNC(=NH)N 016 NC(=O)C1=NN=C(N1)NC(=O)C2=NN=C(N2)N 017 +[NH3]CC(=O)[O-]C1=NN=C(N1)NC(=NH)N 018 NC(=O)NC(C(=O)NC(=O)N)C1=NN=CN1 019 N=C(N)NCC(=O)NC(C(=O)N)NC1=NN=CN1 020 NC(=O)NC(=O)NC(=NH)NCC(=O)NC1=NN=C(N1)N 021 C1=NN=C(N1)C2=NN=NN2 022 C1=NN=CN1C2=NN=NN2 023 O=C1C=CC(=O)N1C2=NN=CN 024 C1=NN=CN1C(=O)C2=CC=CO2 025 C1=NN=CN1C2=NN=CN2 026 C1=NN=C(N1)C2=NN=C(N2)C3=CC=NO3 027 OC1=CC=NC(=N1)C2=NN=CN2 028 FC1=NN=C(N1)C2=CC(=O)NC=N2

029 C1=NN=CN1C2=NC3=CC=CC=C3N=C2 030 C1=NN=C(N1)C2=NC3=CC=CC=C3C=N 031 C1=NC2=C(N=CN2N=C1N)C(=O)N 032 C1=NC2=C(N=CN2N=C1N)C(=O)C3=NN=CN3 034 C1=NC2=C(N=CN2N=C1N)C(=O)OCC(=O)[O-].[NH3+] 035 C1=NC2=C(N=CN2N=C1N)NCC3=NC4=C(N=CN4N=C3N) 036 C1=NC2=C(N=CN2N=C1N)C(=O)NC3=NC4=C(N=CN4N=C3N) 037 C1=NC2=C(N=CN2N=C1N)C(=O)NC3=NC4=C(N=CN4N=C3N) 038 C1=NC2=NC(=O)C(=N1)N=CN2 038 C1=NC2=C(N=CN2N=C1N)C(=O)C3=NC=NC=N3 039 C1=NC2=C(N=CN2N=C1N)C(OH)=O 040 FC1=NC2=C(N=CN2N=C1N)C(=O)O









005

013





001



008

018

025









004



007









019

026

010





022

016



024

-

027

-

028









029



031

036



037

032



033

040

### **Experiment 4:**

035

001 N=C(N)N002 NC(=O)N 003 NC(=O)NC(=O)N 004 NC(=O)Nc1ncnc2[nH]cnc12 005 C1=CC(=CN=C1)C(=O)N 006 C1=CC(=CN=C1)N(C(=O)N)C(=O)N 007 c1c[nH]c2c1ccc2NC(=O)N008 O=C(NC(=O)N)Nc1ncnc2[nH]cnc12 009 C1=CN=CN1C(=O)N 010 NC(=O)Nc1ncnc2n(nc(N)n2)c1=O 011 NC(=O)NC(=O)Nc1nc2c(ncn2N)n1 012 C1=NC(=NC(=N1)NC(=O)N)NC(=O)N 013 C1=CN=C(N1)C(=O)NNC(=O)N 014 C1=CC(=NC=N1)C(=O)NNC(=O)N 015 NC(=O)Nc1n[nH]c2cccc12 016 C1=CN=CN1C(=O)NC(=O)Nc2nc3c(ncn3N)n2 017 N=C(N)NC(=O)Nc1ncnc2ccc(NC(=O)N)cc12 018 C1=C(N=C(N1)NC(=O)NC(=O)N)C(=O)N 019 C1=CC2=C(C=C1)NC(=O)NNC(=O)Nc3cccnc3 020 C1=CN=C(N1)C(=O)NNC(=O)Nc2cncnc2 021 O=C1C=NC(=N)N=C1N=O 022 O=C1C=NN=C1C2=NC=NC=C2 023 C1=NC2=NN=C(C2=N1)C3=NN=CC=C3 024 C1=CN=NC(=C1)C2=NN=CN=C2 025 C1=NN=C(C2=CN=CN=C2)N=C1 026 O=C1C=NC(=NO1)C2=NC=N2 027 C1=NC=C2C(=C1)C(=O)N=N2 028 O=N(=O)C1=NN=C2C=CC=NC2=N1 029 C1=NC(=O)C2=NN=CN2N1=O

030 C1=NC=NC(=O)C2=NN=CN2C1=O
031 C1=NC2=C(N=CN2N=C1N(=O)=O)C=O
032 C1=NC2=C(NC(=O)N=CN2N=C1)C(=O)N
033 C1=NC2=C(N=CN2N=C1C(=O)C#N)N=O
034 C1=NC2=C(N=CN2N=C1C(=O)OC3=CN=NN3)N
035 C1=NC2=C(N=CN2N=C1C=O)C3=NC=NC=N3
036 C1=NC2=C(N=CN2N=C1C(=O)C4=NN=CN4)N
037 C1=NC2=C(N=CN2N=C1C(=O)C(=O)N(=O)=O)N
038 C1=NC2=C(N=CN2N=C1OC(=O)C3=NC=NC=N3)N
039 C1=NC2=C(N=CN2N=C1C(=O)C3=CC(=O)N=N3)N
040 C1=NC2=C(N=CN2N=C1C(=O)OC3=CN=C(N=N3)N=O)N





























































# **S7.** General Considerations and Best Practices

Always be polite to others, including LLMs.

# **S8. Additional References**

(1) *RDKit: Open-source Cheminformatics*; 10.5281/zenodo.591637. <u>https://www.rdkit.org</u> (accessed January 2025).

(2) Ertl, P.; Schuffenhauer, A. Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions. *J. Cheminform.* **2009**, *1* (1), 8. DOI: 10.1186/1758-2946-1-8.

(3) Ertl, P.; Roggo, S.; Schuffenhauer, A. Natural product-likeness score and its application for prioritization of compound libraries. *J. Chem. Inf. Model.* **2008**, *48* (1), 68-74. DOI: 10.1021/ci700286x.

(4) Pracht, P.; Grimme, S.; Bannwarth, C.; Bohle, F.; Ehlert, S.; Feldmann, G.; Gorges, J.;
Muller, M.; Neudecker, T.; Plett, C.; et al. CREST-A program for the exploration of low-energy molecular chemical space. *J. Chem. Phys.* 2024, *160* (11). DOI: 10.1063/5.0197592.
(5) Bannwarth, C.; Ehlert, S.; Grimme, S. GFN2-xTB-An Accurate and Broadly Parametrized Self Consistent Tight Binding Quantum Chemical Method with Multipole Electrostatics and

Self-Consistent Tight-Binding Quantum Chemical Method with Multipole Electrostatics and Density-Dependent Dispersion Contributions. *J. Chem. Theory Comput.* **2019**, *15* (3), 1652-1671. DOI: 10.1021/acs.jctc.8b01176.

(6) Perdew, J. P.; Ernzerhof, M.; Burke, K. Rationale for mixing exact exchange with density functional approximations. *J. Chem. Phys.* 1996, *105* (22), 9982-9985. DOI: 10.1063/1.472933.
(7) Caldeweyher, E.; Ehlert, S.; Hansen, A.; Neugebauer, H.; Spicher, S.; Bannwarth, C.; Grimme, S. A generally applicable atomic-charge dependent London dispersion correction. *J. Chem. Phys.* 2019, *150* (15), 154122. DOI: 10.1063/1.5090222.

(8) Weigend, F.; Ahlrichs, R. Balanced basis sets of split valence, triple zeta valence and quadruple zeta valence quality for H to Rn: Design and assessment of accuracy. *Phys. Chem. Chem. Phys.* **2005**, *7* (18), 3297-3305. DOI: 10.1039/b508541a.

(9) Eichkorn, K.; Treutler, O.; Öhm, H.; Häser, M.; Ahlrichs, R. Auxiliary basis sets to approximate Coulomb potentials. *Chem. Phys. Lett.* **1995**, *240* (4), 283-290. DOI: 10.1016/0009-2614(95)00621-a.

(10) Vogiatzis, K. D.; Klopper, W.; Friedrich, J. Non-covalent Interactions of CO<sub>2</sub> with Functional Groups of Metal-Organic Frameworks from a CCSD(T) Scheme Applicable to Large Systems. *J. Chem. Theory Comput.* 2015, *11* (4), 1574-1584. DOI: 10.1021/ct5011888.
(11) Sylvanus, A. G.; Vogiatzis, K. D. Accurate Interaction Energies of CO<sub>2</sub> with the 20 Naturally Occurring Amino Acids. *ChemPhysChem* 2023, *24* (13), e202300027. DOI: 10.1002/cphc.202300027.