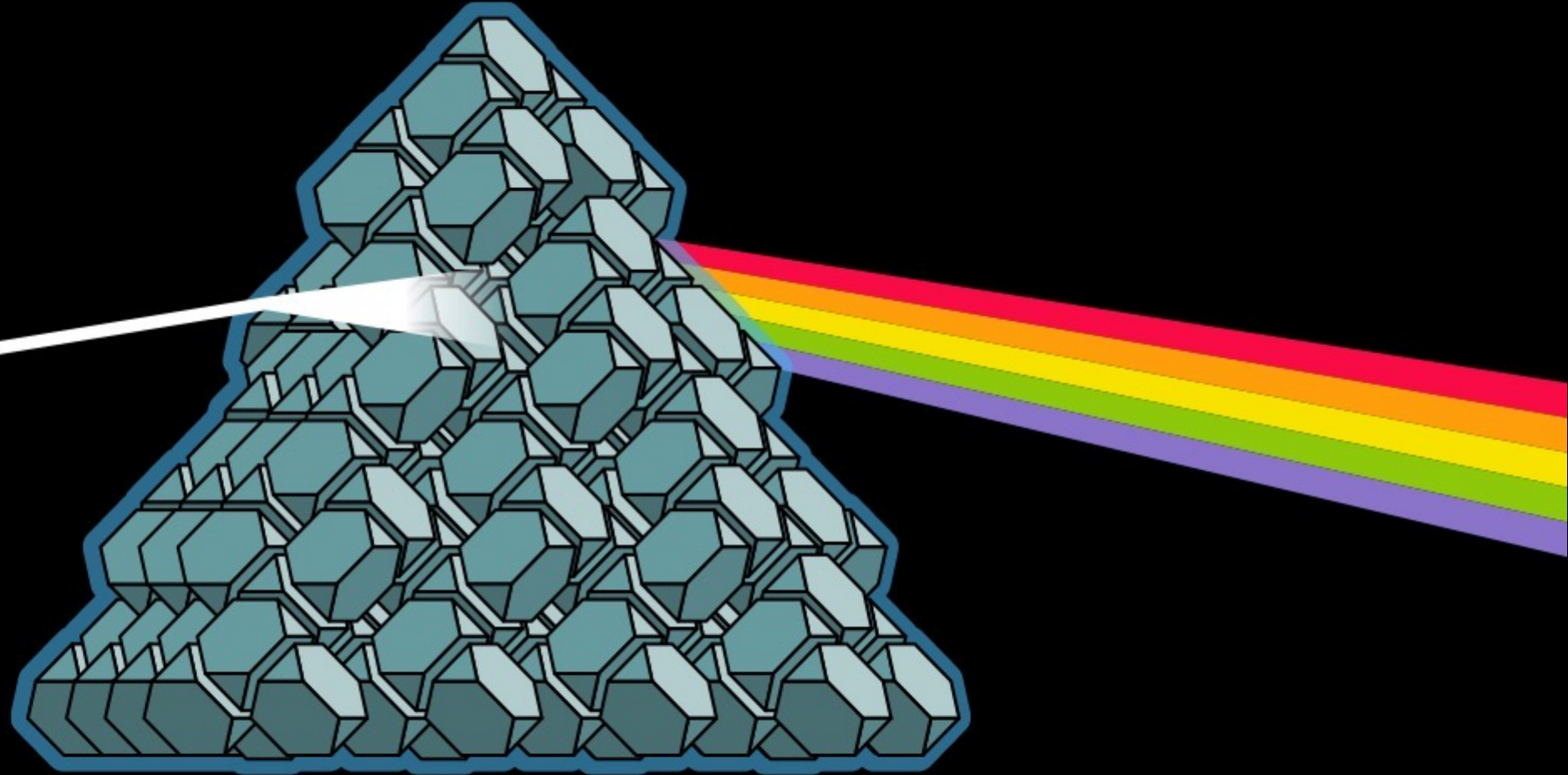# The Search for Novel Mesoscale Materials
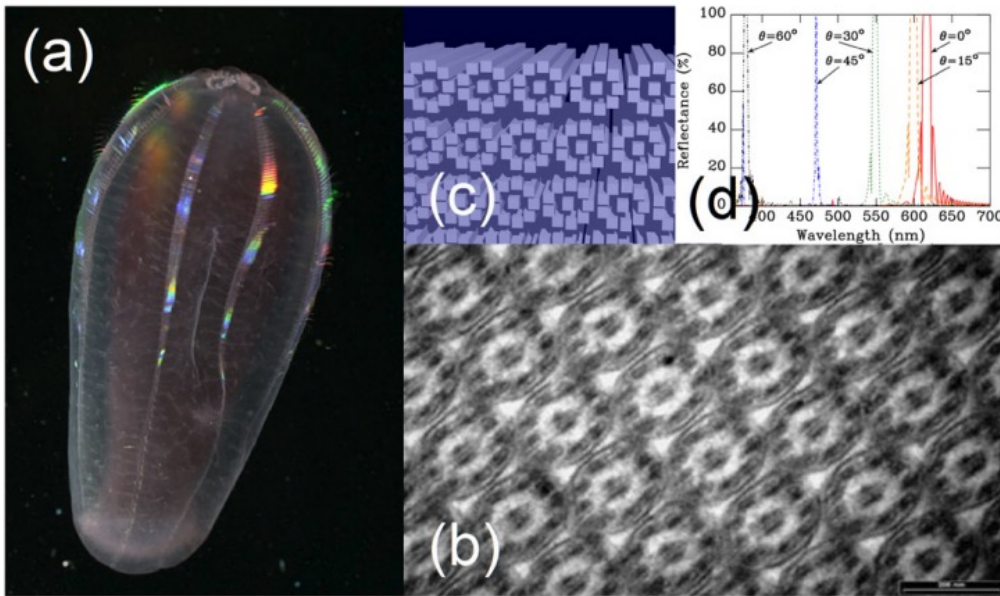
**Rose K. Cersonsky**

Laboratory of Computational Science and Modeling (COSMO)

École Polytechnique Fédérale de Lausanne (EPFL)
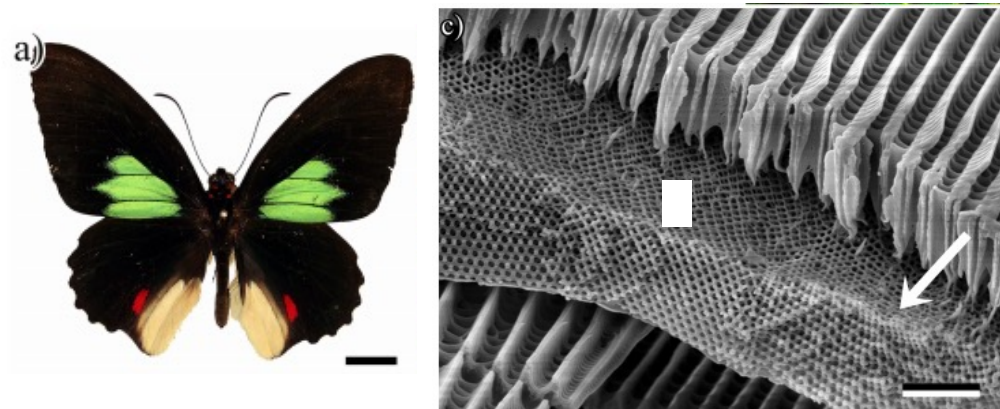
Lausanne, Switzerland

Designing Nanoparticles for the Self-Assembly of Novel (Photonic) Materials

**Optical properties of the iridescent organ of the comb-jellyfish Beroë cucumis (Ctenophora)**
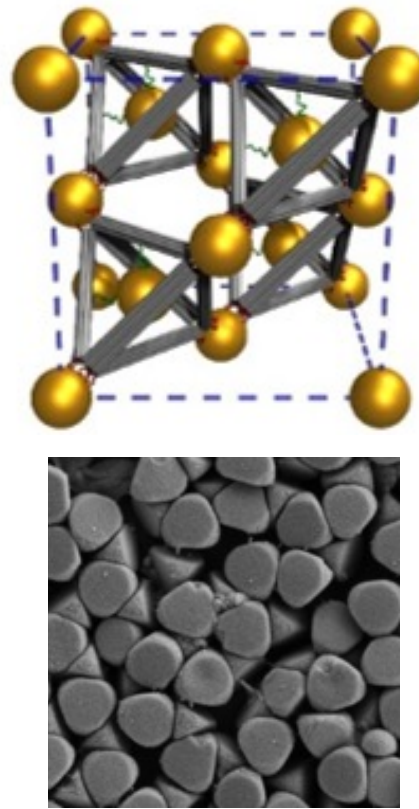Victoria Welch, *et al.*
Phys. Rev. E 73, 041916 2006

**Optical properties of gyroid structured materials: from photonic crystals to metamaterials**
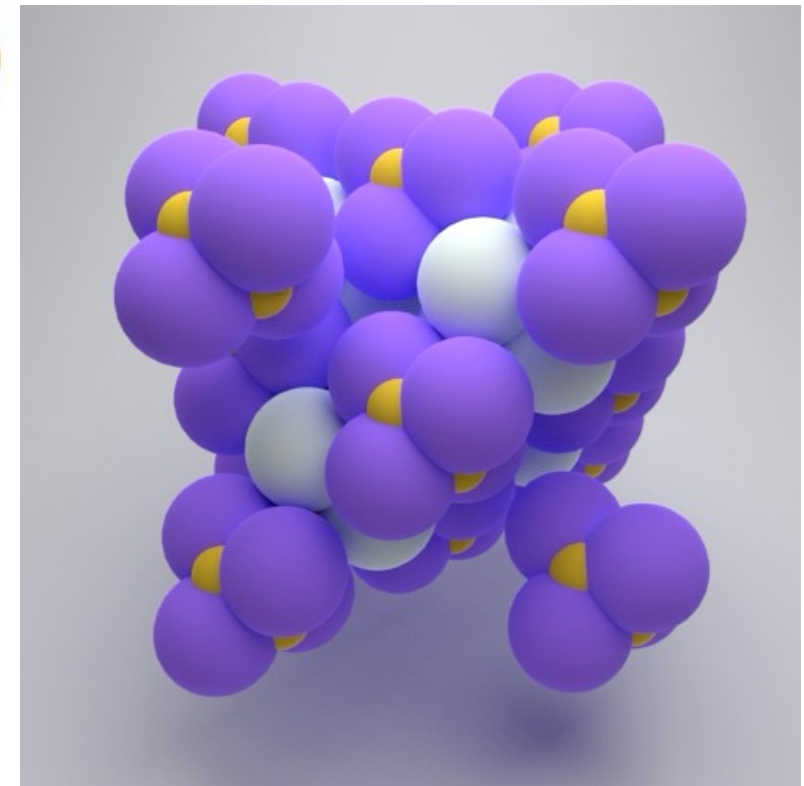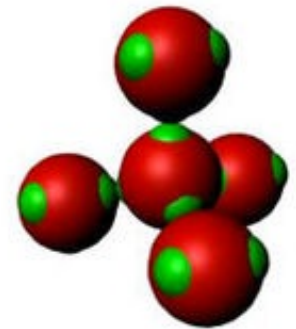James A. Dolan , *et al.*
Advanced Optical Materials 3 (1), 12-32

Colloidal crystals with diamond symmetry at optical lengthscales
Yifan Wang, et al.
Nature Comm. 8, 14173 (2017)

Entropy driven assembly of truncated colloidal tetrahedra into diamond structure
Zhe Gong, et al.

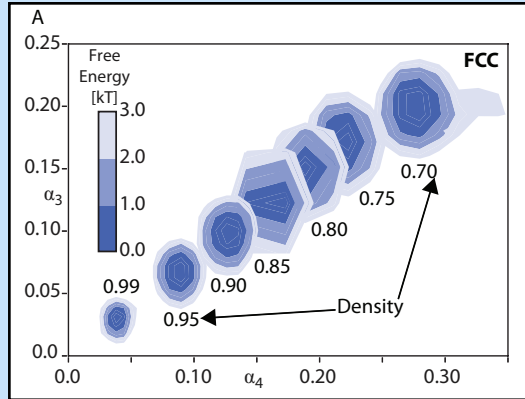Colloidal Diamond
He, M., et al.
Nature 585, 524-529 (2020).

Diamond family of nanoparticle superlattices
W. Liu, et. al,
Science 351, 582-586 (2016).

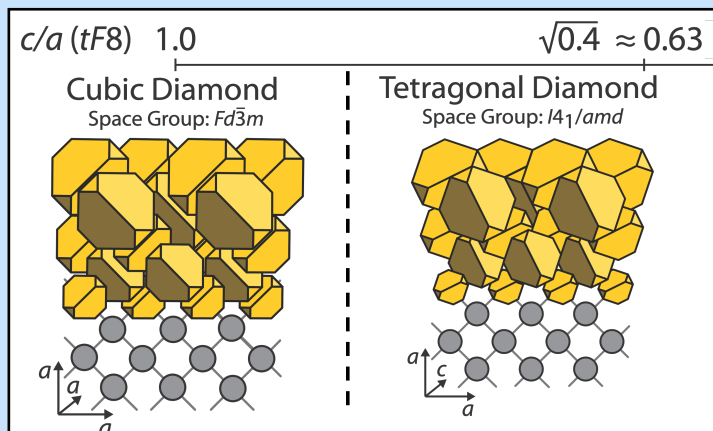## Relevance of packing to colloidal self-assembly.

Cersonsky, R. K., van Anders, G., Dodd, P. M., & Glotzer, S. C. (2018). *Proceedings of the National Academy of Sciences, 115*(7), 1439-1444.



- Pauling's packing rules are **not** a causal mechanism for nanoparticle self-assembly

- Using the Digital Alchemy framework, I showed that **adding small imperfections** to nanoparticle shapes would better stabilize nanocrystals

## Pressure-Tunable Photonic Band Gaps in an Entropic Colloidal Crystal

Cersonsky, R. K., Dshemuchadse, J., Antonaglia, J., van Anders, G., & Glotzer, S. C. (2018). *Physical Review Materials, 2*(12), 125201.



- Nanoparticles that stabilized diamond in self-assembly can **transition to lower-symmetry derivatives** at high pressure

- Small distortions in diamond **did not destroy the photonic band gap**

At time of presentation, this manuscript was not yet published, please see rosecersonsky.com for recent publications.

…Small distortions in diamond **did not destroy the photonic band gap…**
…minimal effect on the photonic band structure…

**what is the span of crystallographic structures capable of supporting a photonic band gap?**



**"Direct"**
*high* dielectric medium on lattice sites

**"Inverse"**
*low* dielectric medium on lattice sites

**Scraped 1300 unique crystal structures from crystal repositories**

**Ignoring atomic species, turned each structure into 2 templates with which to sample multiple parameters**

**Ran >150,000 band structure to determine which "templates" supported PBGs**

351 Photonic "Templates"

474 Unique Gaps

Database of Photonic Crystals:
https://glotzerlab.engin.umich.edu/photonics/index.html

Appendix of Band Structures:
https://deepblue.lib.umich.edu/handle/2027.42/153520

ε = 4
ε = 6
ε = 8
ε = 10
ε = 12
ε = 14
ε = 16

Each circle represents the maximum gap (circle size) found for a given template (radius), dielectric contrast (ring), and band location (color).

Each circle represents the maximum gap (circle size) found for a given template (radius), dielectric contrast (ring), and band location (color).

# Inverse AB₁₃
## Maximum Gaps: 13.3% (Gaps 5-6), 4.78% (Gaps 10-11)

# Inverse Clathrate-II
## Maximum Gap: 33.9%



Inverse AB₁₃ (cP26)

Inverse Clathrate-II (cF136)

Inverse AB₁₃ (cP26)

ε = 4
ε = 6
ε = 8
ε = 10
ε = 12
ε = 14
ε = 16

G    H    Cage A    I

Cage D

Clathrate colloidal crystals.
Lin, H., Lee, S., Sun, L., Spellings, M.,
Engel, M., Glotzer, S. C., & Mirkin, C. A.
Science, 355(6328), 931-935.

Each circle represents the maximum gap (circle size) found for a given template (radius), dielectric contrast (ring), and band location (color).

Lithium Oxide (cF12)

c-OT$_2$

Self-assembly of a space-tessellating structure in the binary system of hard tetrahedra and octahedra.
Cadotte, Andrew T., et al.
Soft matter 12.34 (2016): 7073-7078.

Lithium Oxide (cF12)

$\varepsilon = 4$
$\varepsilon = 6$
$\varepsilon = 8$
$\varepsilon = 10$
$\varepsilon = 12$
$\varepsilon = 14$
$\varepsilon = 16$

The lithium-oxide structure (a.k.a. Fluorite, c-OT$_2$, and F-RD) exhibits photonic anomalies, including a band gap that is largest at lower dielectric contrast.

Each circle represents the maximum gap (circle size) found for a given template (radius), dielectric contrast (ring), and band location (color).

# How else can we use this large dataset?



Database of Photonic Crystals:
https://glotzerlab.engin.umich.edu/photonics/index.html

203: Innovations in Methods of Data Science
Monday, November 8, 2021, 4:30 PM - 4:45 PM
Marriott Copley Place - Salon H/I

At time of presentation, this manuscript was not yet published, please see rosecersonsky.com for recent publications.

At time of presentation, this work was not yet published, please see rosecersonsky.com for recent publications.

# The Search for Novel Mesoscale Materials

**Rose K. Cersonsky**

Laboratory of Computational Science and Modeling (COSMO)
École Polytechnique Fédérale de Lausanne (EPFL)
Lausanne, Switzerland

At time of presentation, this work was not yet published, please see rosecersonsky.com for recent publications.

**Come see my last talk!**
**203e - Improving Data Sub-Selection for Supervised Tasks with Principal Covariates Regression**
Monday, November 8, 2021
4:30 PM - 4:45 PM EDT
Marriott Copley Place - Salon H/I

**rosecersonsky.com**