# Mathematics based graph neural network for drug design
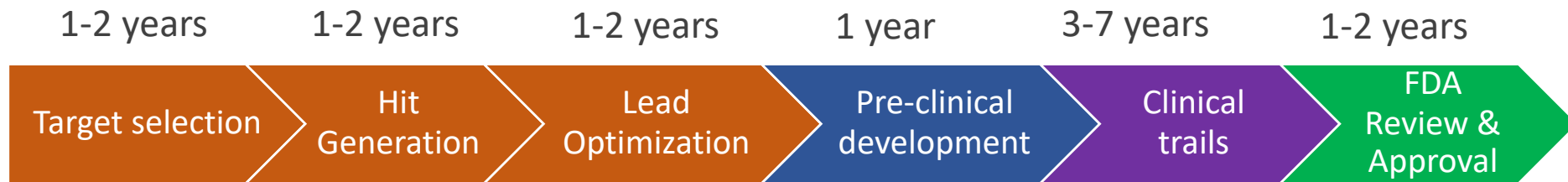
Duc Nguyen

Department of Mathematics

University of Kentucky

The Fourth TSIMF Conference on Computational and Mathematical Bioinformatics and Biophysics
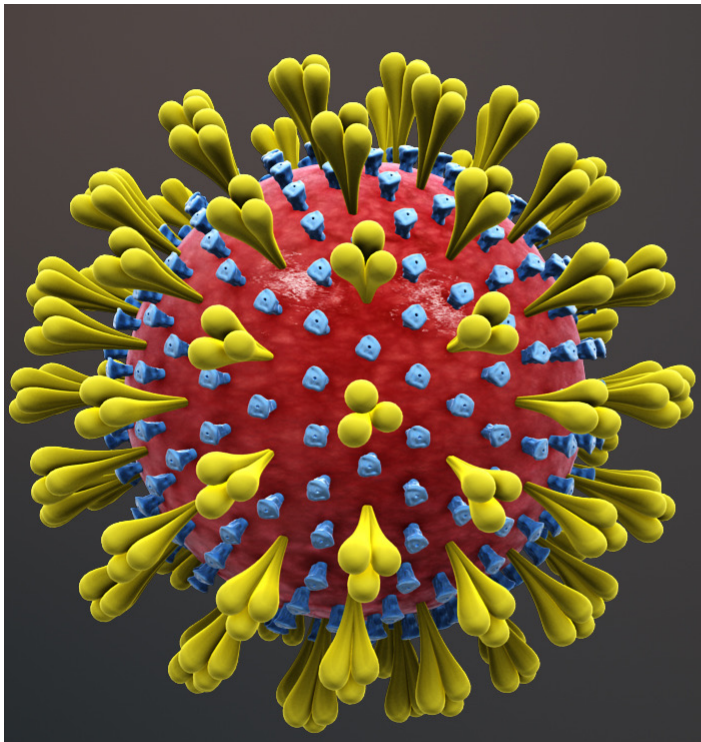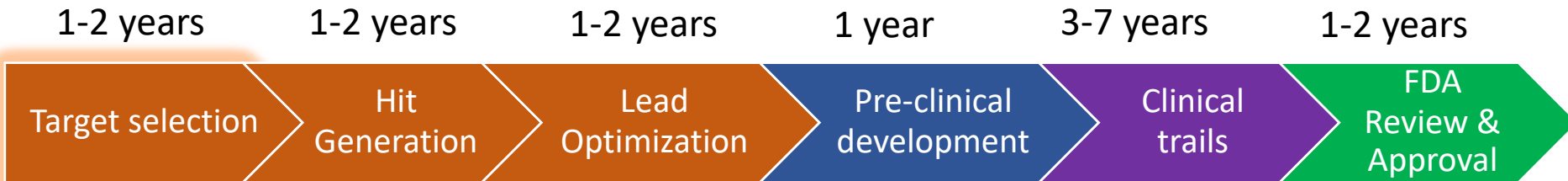
December 12-15, 2021

# Rational Drug Discovery

| 1-2 years | 1-2 years | 1-2 years | 1 year | 3-7 years | 1-2 years |
|---|---|---|---|---|---|
| Target selection | Hit Generation | Lead Optimization | Pre-clinical development | Clinical trails | FDA Review & Approval |

# Rational Drug Discovery

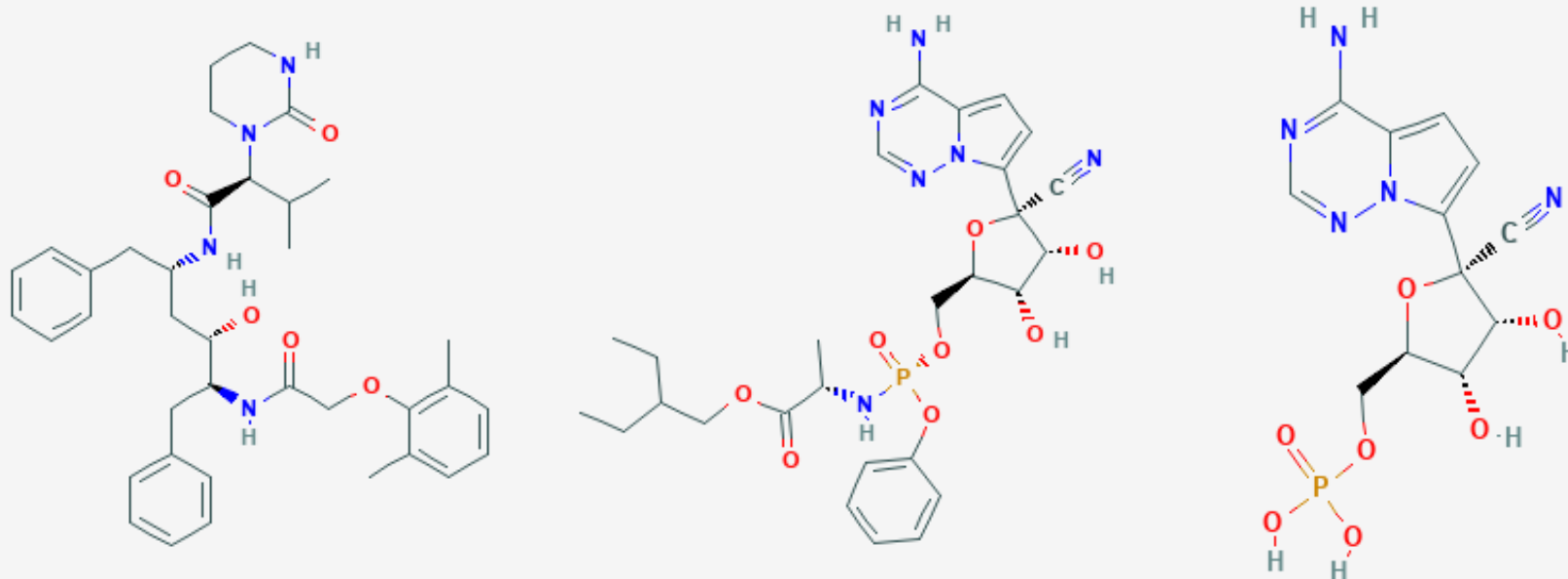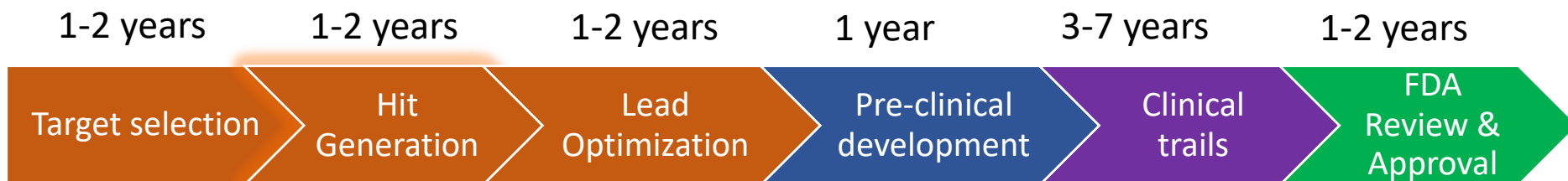| 1-2 years | 1-2 years | 1-2 years | 1 year | 3-7 years | 1-2 years |
|-----------|-----------|-----------|--------|-----------|-----------|
| Target selection | Hit Generation | Lead Optimization | Pre-clinical development | Clinical trails | FDA Review & Approval |



- **COVID-19 (SARS-CoV-2)**

- First reported on Dec 30, 2019

- Global health emergency by WHO on 01-29-20

- As of 12-14-21: 5.31M dead cases, > 271M infected cases

https://en.wikipedia.org/wiki/Coronavirus

# Rational Drug Discovery

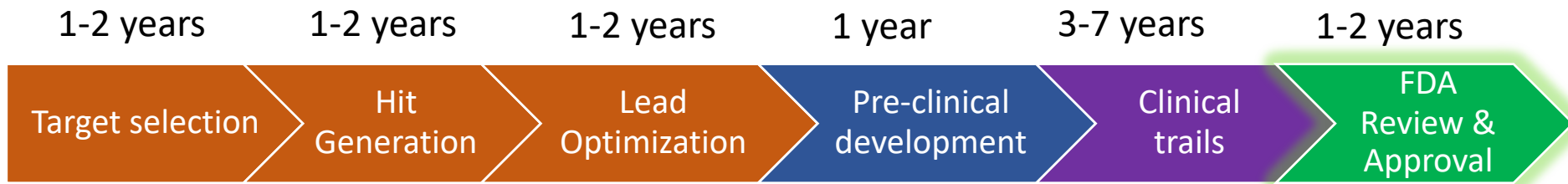| 1-2 years | 1-2 years | 1-2 years | 1 year | 3-7 years | 1-2 years |
|---|---|---|---|---|---|
| Target selection | Hit Generation | Lead Optimization | Pre-clinical development | Clinical trails | FDA Review & Approval |

- Improve potency
- reduced off-target activities
- Reasonable physiochemical/metabolic properties *in vivo* pharmacokinetics

# Rational Drug Discovery

| 1-2 years | 1-2 years | 1-2 years | 1 year | 3-7 years | 1-2 years |
|---|---|---|---|---|---|
| Target selection | Hit Generation | Lead Optimization | Pre-clinical development | Clinical trails | FDA Review & Approval |



Important properties: binding affinity (IC50), toxicity, solubility, …

# Rational Drug Discovery

| 1-2 years | 1-2 years | 1-2 years | 1 year | 3-7 years | 1-2 years |
|---|---|---|---|---|---|
| Target selection | Hit Generation | Lead Optimization | Pre-clinical development | Clinical trails | FDA Review & Approval |

**?**

No available FDA-approved drugs for COVID-19 on market yet!

# Rational Drug Discovery

| 1-2 years | 1-2 years | 1-2 years | 1 year | 3-7 years | 1-2 years |
|---|---|---|---|---|---|
| Target selection | Hit Generation | Lead Optimization | Pre-clinical development | Clinical trails | FDA Review & Approval |

❌ Lengthy process ( > 10 years)

❌ Expensive ( > $2.6 billion)

❌ High failure rate

# Rational Drug Discovery

| 1-2 years | 1-2 years | 1-2 years | 1 year | 3-7 years | 1-2 years |
|-----------|-----------|-----------|--------|-----------|-----------|
| Target selection | Hit Generation | Lead Optimization | Pre-clinical development | Clinical trails | FDA Review & Approval |

- ✕ Lengthy process ( > 10 years)

- ✕ Expensive ( > $2.6 billion)

- ✕ High failure rate

Great opportunities for Math and AI

Protein Data Bank, as of 2021: 174,994 structures

230 million compounds — ZINC

1.9 million compounds — ChEMBL

# Molecular Representations

| SMILES String | 2D Diagram | 3D |
|---|---|---|

CC(F)(F)[C@H]1OC(N)=
N[C@](C)(c2cc(NC(=O)c
3cnc(OCF)cn3)ccc2F)[C
@H]1F

$\mathcal{F}_1$

$\mathcal{F}_2$

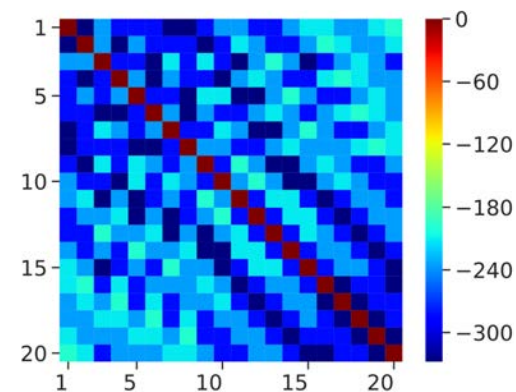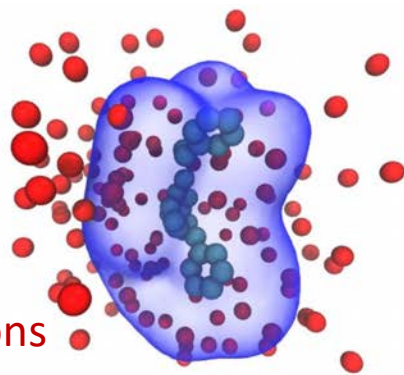$\mathcal{F}_3$

Fingerprint

Fingerprint

Fingerprint

# Molecular Representations

(Nguyen, Cang, Wei, PCCP 2020, …)

3D

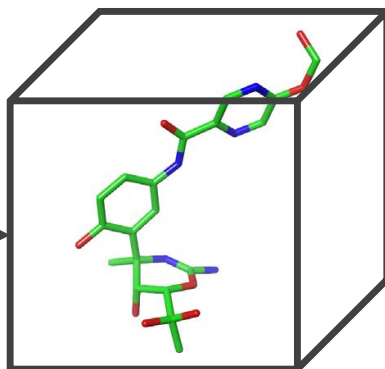Mathematical Representations

(Meng and Xia, Sci. Adv. 2021, …)
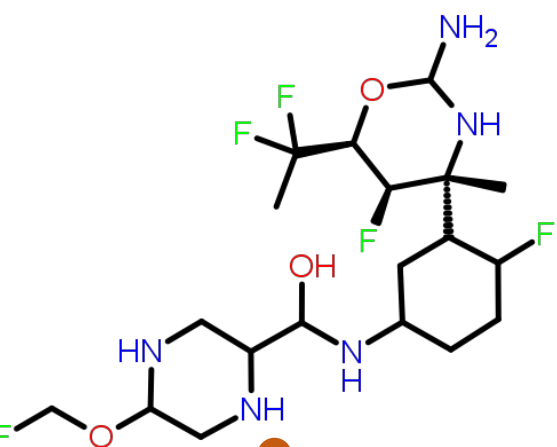
$\beta_0$

$\beta_1$

$\beta_2$

3D-Image Like Representations

Jimenez et. al. (2018),
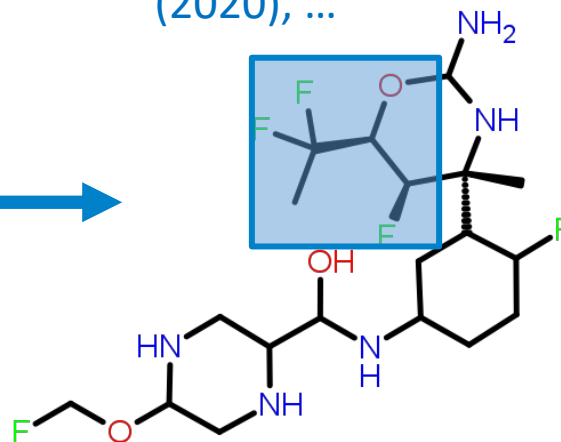Hassan-Harrirou et. al. (2020), …

# Molecular Representations

## 2D Diagram



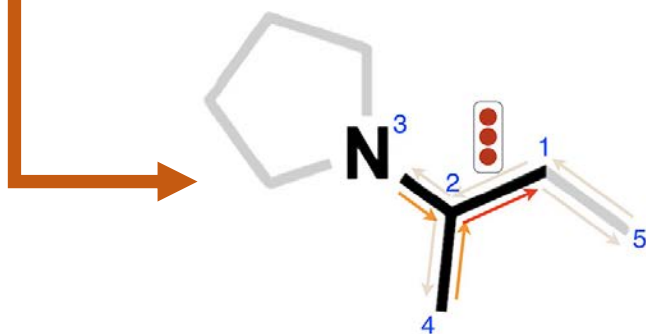2D Image

Goh et. al. (2018), Rifaioglu et. al. (2020), …

Graph Neural Network

Duvenaud et. al. (2015), Kearnes et. al. (2016), Wu el. al. (2018), Yang el. al. (2019), Withnall (2020), …

# Molecular Representations

## SMILES String

CC(F)(F)[C@H]1OC(N)=
N[C@](C)(c2cc(NC(=O)c
3cnc(OCF)cn3)ccc2F)[C
@H]1F

Grow et. al., CIS 2019,
Gao et. al., JCIM 2020

**GRU**  **LS**  **GRU**

Deep Learning's
Hidden States

(Dong Chen et. al., Nat. Com. 2021)

BT-FP

Fine-turned model

<s> N = C C N </s>

N = C C N

Task-specific data

2D Fingerprints

Language Embedding
(Word2vec, …)

FP2, Daylight, MACCS,
Estate, ECPF4,
Pharm2D, ERG …

(Gao et. al., PCCP 2020)

Rogers and Hahn (2010), …

Jaeger et. al. (2018)
Goh et. al. (2018), …

# Mathematical Graphs for Molecules
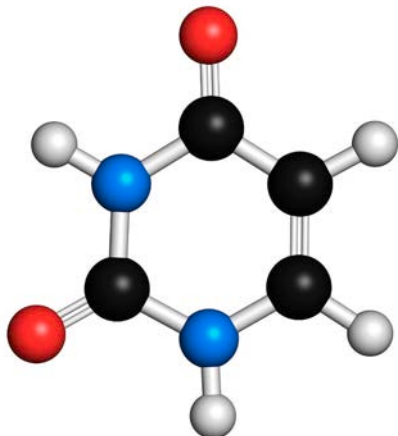
🔴 Oxygen   ⚫ Carbon   🔵 Nitrogen   ⚪ Hydrogen

Uracil

Subgraphs: $\mathcal{G}_{O,C}, \mathcal{G}_{O,N}, \mathcal{G}_{O,H},$
$\mathcal{G}_{C,N}, \mathcal{G}_{C,H}$

$$A = \begin{bmatrix} 0 & w_{12} & 0 & w_{14} \\ w_{21} & 0 & w_{23} & 0 \\ 0 & w_{32} & 0 & w_{34} \\ w_{41} & 0 & w_{42} & 0 \end{bmatrix}$$

$$A = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{bmatrix}$$

Weighted Adjacency Matrix

Adjacency Matrix

$\mathcal{G}_{O,N}$

- $w_{ij}(d) = e^{-\left(\frac{d}{\eta}\right)^{k}}$

- $w_{ij}(d) = \dfrac{1}{1 + \left(\frac{d}{\eta}\right)^{v}}$

$$L = D - A = \begin{bmatrix} 2 & -1 & 0 & -1 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ -1 & 0 & -1 & 2 \end{bmatrix}$$

Laplacian Matrix

(Nguyen, Wei, JCIM 2019)

# Persistent Spectral Graph: Graph + Topology

- Simplexes:



**0-simplex**    **1-simplex**    **2-simplex**    **3-simplex**
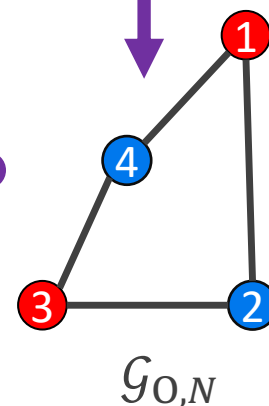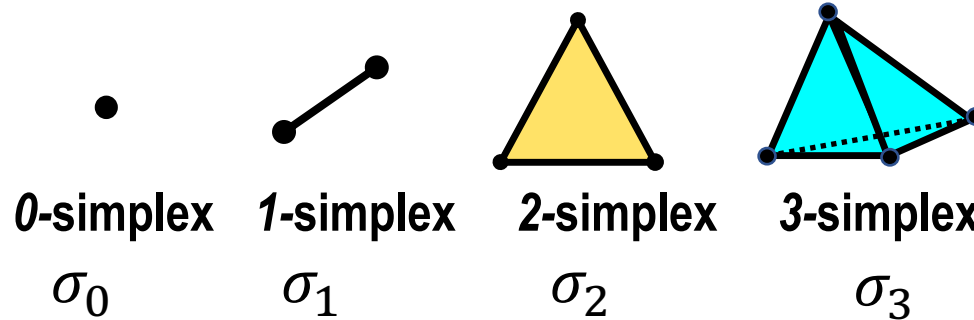
$\sigma_0$      $\sigma_1$      $\sigma_2$      $\sigma_3$

- $q$-chain: $\quad \displaystyle\sum_j w_j \sigma_q^j, \quad w_j \in \mathbb{Z}_2, \sigma_q^j \in K$

Wang, Nguyen, Wei (IJNMBE, 2020)

- $q$-chain Group: $C_q(K) = \left( \left\{ \displaystyle\sum_j w_j \sigma_q^j \right\}, + \right)$

Meng and Xia, (Sci. Adv. 2021)

- Boundary operator:

$$\partial_q : C_q(K) \to C_{q-1}(K) \quad , \partial_q \, \sigma_q = \sum_{j=0}^{q} (-1)^j \langle v_0, v_1, \ldots, \widehat{v_j}, \ldots, v_q \rangle$$

- Adjoint boundary operator: $\quad \partial_q^* : C_{q-1}(K) \to C_q(K)$

- $q$-combinatorial Laplacian operator: $\quad \Delta_q = \partial_{q+1} \partial_{q+1}^* + \partial_q^* \partial_q$

- $q$-combinatorial Laplacian matrix: $\quad \mathcal{L}_q = \mathcal{B}_{q+1} \mathcal{B}_{q+1}^T + \mathcal{B}_q^T \mathcal{B}_q$

- Betti numbers: $\quad \beta_q =$ # of zeros eigenvalues of $\mathcal{L}_q(K)$

## Filtration on fullerene $C_{20}$



### Accumulated combinatorial Laplacian matrix for $C_{20}$



## Barcodes of protein-ligand

PDBID: 1hmk



$\beta_0$

$\beta_1$

$\beta_2$

# Differential Geometry

- Element interactive density: $\rho_{kk'}(\mathbf{r}, \eta_{kk'}) = \sum_j w_j \Phi(\|\mathbf{r} - \mathbf{r}_j\|; \eta_{kk'})$, with $\|\mathbf{r} - \mathbf{r}_j\| > r_i + r_j + \sigma$

- Element interactive manifolds (EIMs): $\rho_{kk'}(r, \eta_{kk'}) = c\rho_{\max}, 0 \leq c \leq 1$

- Element interactive curvatures, element interactive areas, ...

## Element Interactive Manifolds

# Mathematics based Deep Learning Models (MathDL)



Protein-ligand complex → Element specific groups → Various Mathematical features → Machine learning prediction

Algebraic topology and/or Differential geometry and/or Graph Theory

(Nguyen et. al., JCAMD 2019)    (Patent No.: US 2019 / 0304568 A1)

# Performance of MathDL in Virtual Screening, Docking, Affinity Ranking

**Red Color** is our Model

(Nguyen, Wei, JCIM 2019)



a) Scoring Power (CASF-2013) — Pearson's R

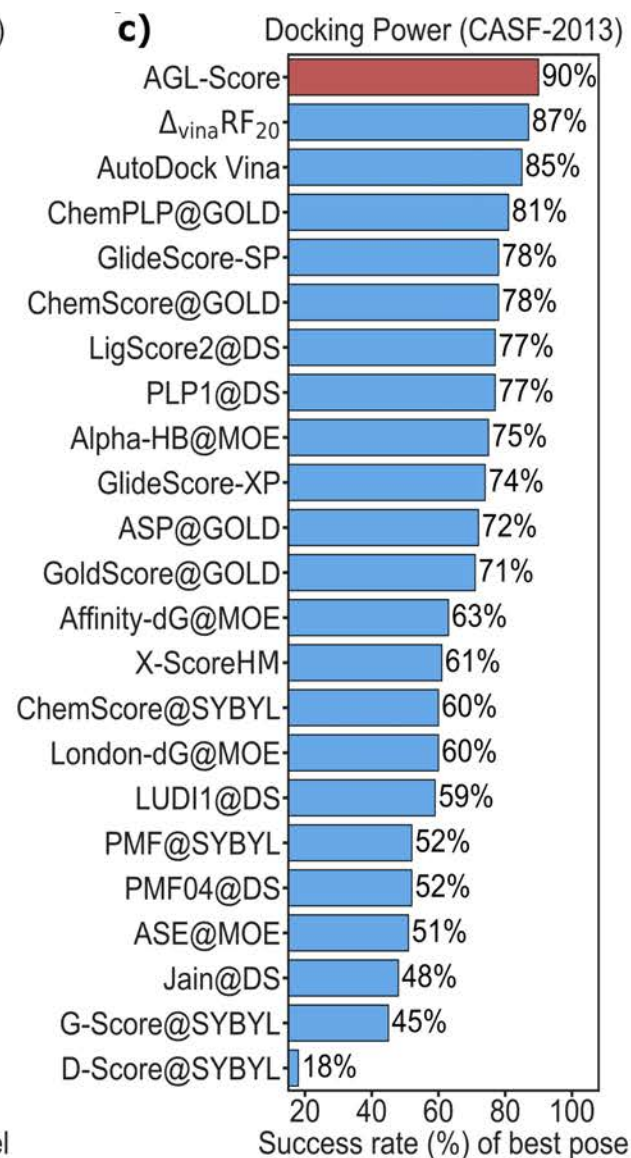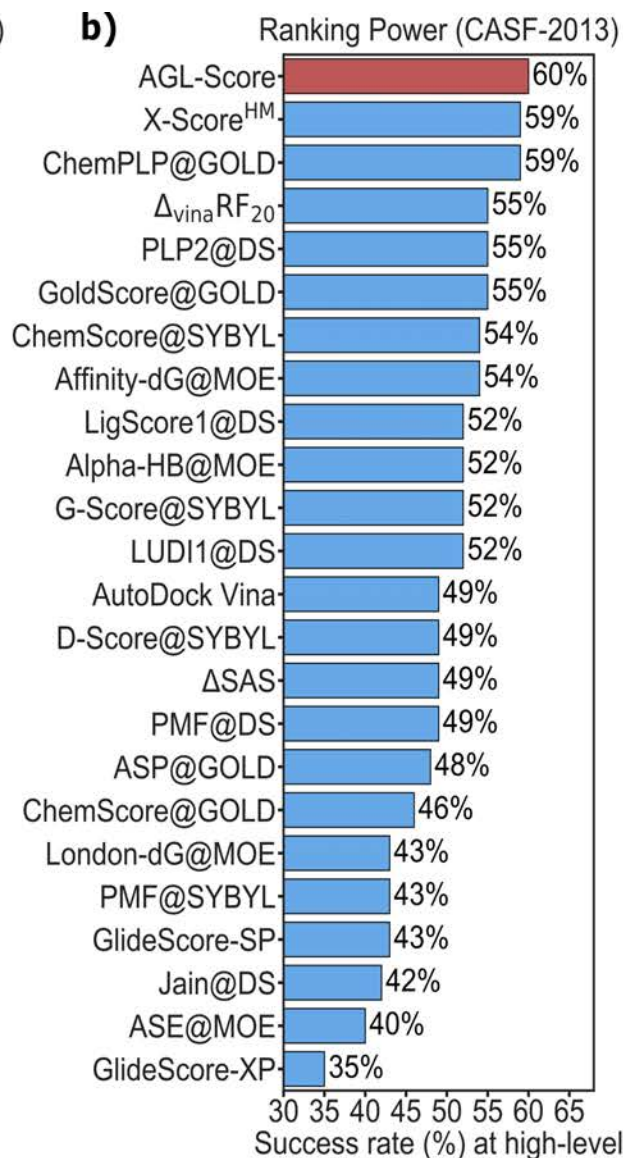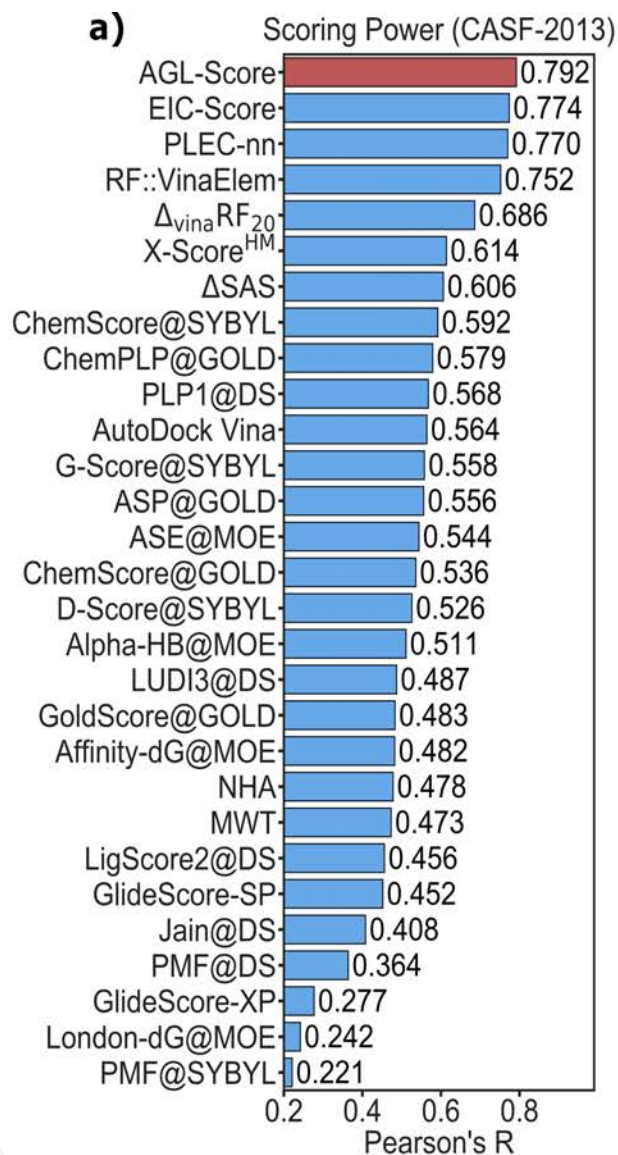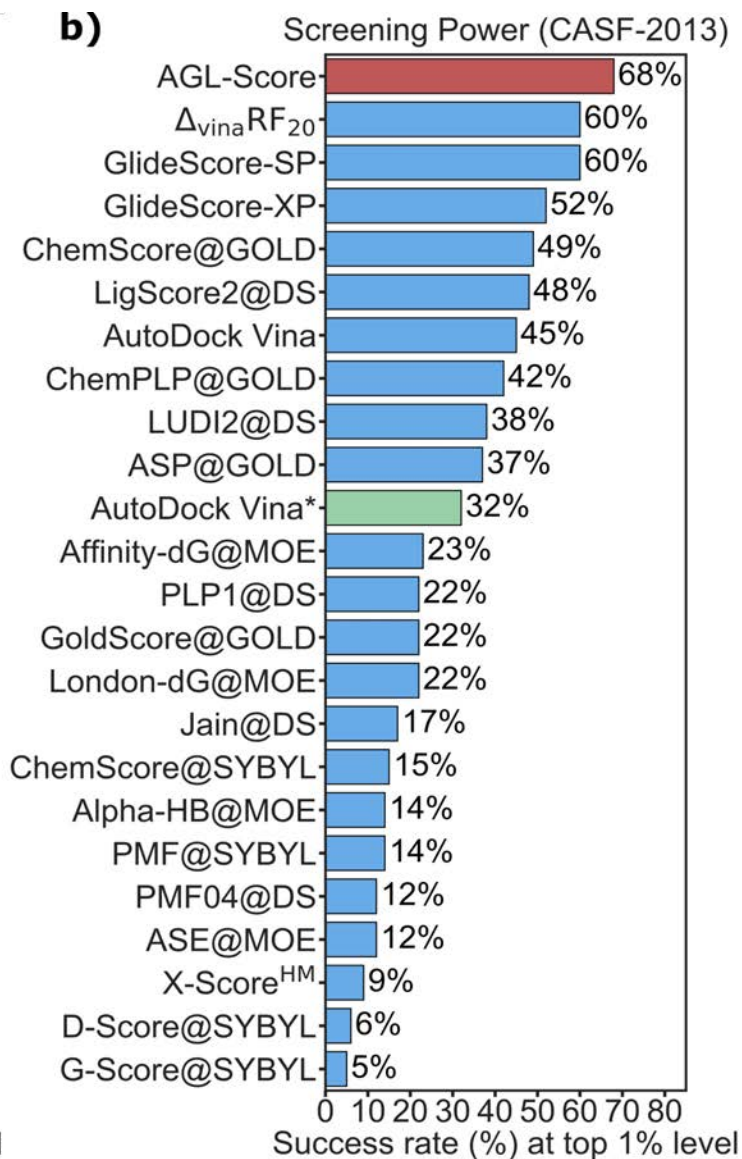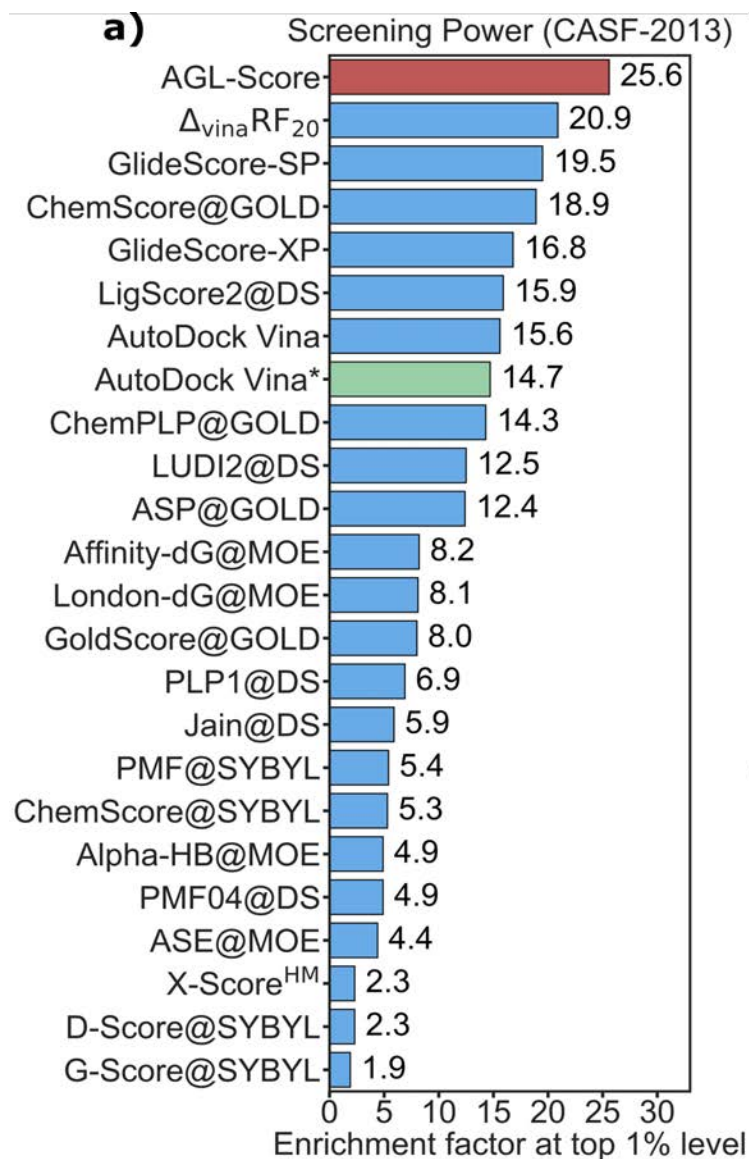| Method | Value |
|---|---|
| AGL-Score | 0.792 |
| EIC-Score | 0.774 |
| PLEC-nn | 0.770 |
| RF::VinaElem | 0.752 |
| $\Delta_{vina}RF_{20}$ | 0.686 |
| X-Score$^{HM}$ | 0.614 |
| $\Delta$SAS | 0.606 |
| ChemScore@SYBYL | 0.592 |
| ChemPLP@GOLD | 0.579 |
| PLP1@DS | 0.568 |
| AutoDock Vina | 0.564 |
| G-Score@SYBYL | 0.558 |
| ASP@GOLD | 0.556 |
| ASE@MOE | 0.544 |
| ChemScore@GOLD | 0.536 |
| D-Score@SYBYL | 0.526 |
| Alpha-HB@MOE | 0.511 |
| LUDI3@DS | 0.487 |
| GoldScore@GOLD | 0.483 |
| Affinity-dG@MOE | 0.482 |
| NHA | 0.478 |
| MWT | 0.473 |
| LigScore2@DS | 0.456 |
| GlideScore-SP | 0.452 |
| Jain@DS | 0.408 |
| PMF@DS | 0.364 |
| GlideScore-XP | 0.277 |
| London-dG@MOE | 0.242 |
| PMF@SYBYL | 0.221 |

b) Ranking Power (CASF-2013) — Success rate (%) at high-level

| Method | Value |
|---|---|
| AGL-Score | 60% |
| X-Score$^{HM}$ | 59% |
| ChemPLP@GOLD | 59% |
| $\Delta_{vina}RF_{20}$ | 55% |
| PLP2@DS | 55% |
| GoldScore@GOLD | 55% |
| ChemScore@SYBYL | 54% |
| Affinity-dG@MOE | 54% |
| LigScore1@DS | 52% |
| Alpha-HB@MOE | 52% |
| G-Score@SYBYL | 52% |
| LUDI1@DS | 52% |
| AutoDock Vina | 49% |
| D-Score@SYBYL | 49% |
| $\Delta$SAS | 49% |
| PMF@DS | 49% |
| ASP@GOLD | 48% |
| ChemScore@GOLD | 46% |
| London-dG@MOE | 43% |
| PMF@SYBYL | 43% |
| GlideScore-SP | 43% |
| Jain@DS | 42% |
| ASE@MOE | 40% |
| GlideScore-XP | 35% |

c) Docking Power (CASF-2013) — Success rate (%) of best pose

| Method | Value |
|---|---|
| AGL-Score | 90% |
| $\Delta_{vina}RF_{20}$ | 87% |
| AutoDock Vina | 85% |
| ChemPLP@GOLD | 81% |
| GlideScore-SP | 78% |
| ChemScore@GOLD | 78% |
| LigScore2@DS | 77% |
| PLP1@DS | 77% |
| Alpha-HB@MOE | 75% |
| GlideScore-XP | 74% |
| ASP@GOLD | 72% |
| GoldScore@GOLD | 71% |
| Affinity-dG@MOE | 63% |
| X-ScoreHM | 61% |
| ChemScore@SYBYL | 60% |
| London-dG@MOE | 60% |
| LUDI1@DS | 59% |
| PMF@SYBYL | 52% |
| PMF04@DS | 52% |
| ASE@MOE | 51% |
| Jain@DS | 48% |
| G-Score@SYBYL | 45% |
| D-Score@SYBYL | 18% |

# Performance of MathDL in Virtual Screening, Docking, Affinity Ranking

**Red Color** is our Model

(Nguyen, Wei, JCIM 2019)



**a)** Screening Power (CASF-2013)

| Method | Enrichment factor at top 1% level |
|---|---|
| AGL-Score | 25.6 |
| $\Delta_{vina}RF_{20}$ | 20.9 |
| GlideScore-SP | 19.5 |
| ChemScore@GOLD | 18.9 |
| GlideScore-XP | 16.8 |
| LigScore2@DS | 15.9 |
| AutoDock Vina | 15.6 |
| AutoDock Vina* | 14.7 |
| ChemPLP@GOLD | 14.3 |
| LUDI2@DS | 12.5 |
| ASP@GOLD | 12.4 |
| Affinity-dG@MOE | 8.2 |
| London-dG@MOE | 8.1 |
| GoldScore@GOLD | 8.0 |
| PLP1@DS | 6.9 |
| Jain@DS | 5.9 |
| PMF@SYBYL | 5.4 |
| ChemScore@SYBYL | 5.3 |
| Alpha-HB@MOE | 4.9 |
| PMF04@DS | 4.9 |
| ASE@MOE | 4.4 |
| X-Score[HM] | 2.3 |
| D-Score@SYBYL | 2.3 |
| G-Score@SYBYL | 1.9 |

**b)** Screening Power (CASF-2013)

| Method | Success rate (%) at top 1% level |
|---|---|
| AGL-Score | 68% |
| $\Delta_{vina}RF_{20}$ | 60% |
| GlideScore-SP | 60% |
| GlideScore-XP | 52% |
| ChemScore@GOLD | 49% |
| LigScore2@DS | 48% |
| AutoDock Vina | 45% |
| ChemPLP@GOLD | 42% |
| LUDI2@DS | 38% |
| ASP@GOLD | 37% |
| AutoDock Vina* | 32% |
| Affinity-dG@MOE | 23% |
| PLP1@DS | 22% |
| GoldScore@GOLD | 22% |
| London-dG@MOE | 22% |
| Jain@DS | 17% |
| ChemScore@SYBYL | 15% |
| Alpha-HB@MOE | 14% |
| PMF@SYBYL | 14% |
| PMF04@DS | 12% |
| ASE@MOE | 12% |
| X-Score[HM] | 9% |
| D-Score@SYBYL | 6% |
| G-Score@SYBYL | 5% |

# Collaboration work with Pfizer



**Binding affinity ranking of 362 compounds (fully blind)**

Red Bars are our models

Pfizer pays **$1 million** annually for this software license

Free software

# Drug Design Data Resource (D3R) Grand Challenges

- Funded in part by National Institute of General Medical Sciences
- Hosted at the University of California, San Diego
- Annually since 2015



**GC4 (2018 -2019): 55 research groups**

# Drug Design Data Resource (D3R) Grand Challenges

# D3R Grand Challenge 2 (2016-2017)

**Stage 1**
Pose Predictions (partials)
Scoring (partials)
Free Energy Set 1 (partials)
Free Energy Set 2 (partials)

**Stage 2**
Scoring (partials)
Free Energy Set 1 (partials)
Free Energy Set 2 (partials)



(Nguyen et. al., JCAMD 2018)

# D3R Grand Challenge 3 (2017-2018)

**Pose Prediction**

| **Cathepsin Stage 1A** | **Cathepsin Stage 1B** |
|---|---|
| Pose Predictions (partials) | Pose Prediction |

**Affinity Rankings excluding Kds > 10 µM**

| **Cathepsin Stage 1** | **Cathepsin Stage 2** | |
|---|---|---|
| Scoring (partials) | Scoring (partials) | |
| Free Energy Set | Free Energy Set | |
| **VEGFR2** | **JAK2 SC2** | **p38-α** |
| Scoring (partials) | Scoring (partials) | Scoring |
| **JAK2 SC3** | **TIE2** 🥇🥈 | **ABL1** |
| Scoring | Scoring 🥇🥈🥉 | Scoring (partials) 🥇🥈🥉 |
| Free Energy Set 🥇🥈🥉 | Free Energy Set 2 🥇🥈🥉 | |

**Active / Inactive Classification**

| **VEGFR2** | **JAK2 SC2** | **p38-α** |
|---|---|---|
| Scoring (partials) | Scoring (partials) | Scoring (partials) |
| **JAK2 SC3** | **TIE2** | **ABL1** |
| Scoring | Scoring (partials) 🥇🥈🥉 | Scoring (partials) |
| Free Energy Set 🥇🥈🥉 | Free Energy Set 1 🥇🥈🥉 | |

**Affinity Rankings for Cocrystalized Ligands**

| **Cathepsin Stage 1** | **Cathepsin Stage 2** 🥇🥈 |
|---|---|
| Scoring (partials) 🥇🥈 | Scoring (partials) |
| Free Energy Set 🥇🥈 | Free Energy Set 🥇🥈🥉 |

Cathepsin S     Kinase: p38-α

(Nguyen et. al., JCAMD 2018)

# D3R Grand Challenge 4 (2018-2019)

**Pose Predictions**

**BACE Stage 1A**
Pose Predictions (Partials) 🥇🥈

**BACE Stage 1B**
Pose Prediction (Partials) 🥈🥉

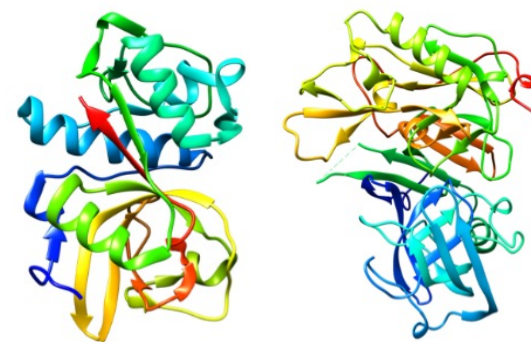**Affinity Predictions**

**Cathepsin Stage 1**
Combined Ligand and Structure Based Scoring 🥇🥈🥉
Ligand Based Scoring (No participation)
Structure Based Scoring 🥇🥈🥉
Free Energy Set 🥇🥈🥉

**BACE Stage 1**
Combined Ligand and Structure (No participation)
Ligand Based Scoring (Partials) (No participation)
Structure Based Scoring (Partials)(No participation)
Free Energy Set (No participation)

**BACE Stage 2**
Combined Ligand and Structure
Ligand Based Scoring (No participation)
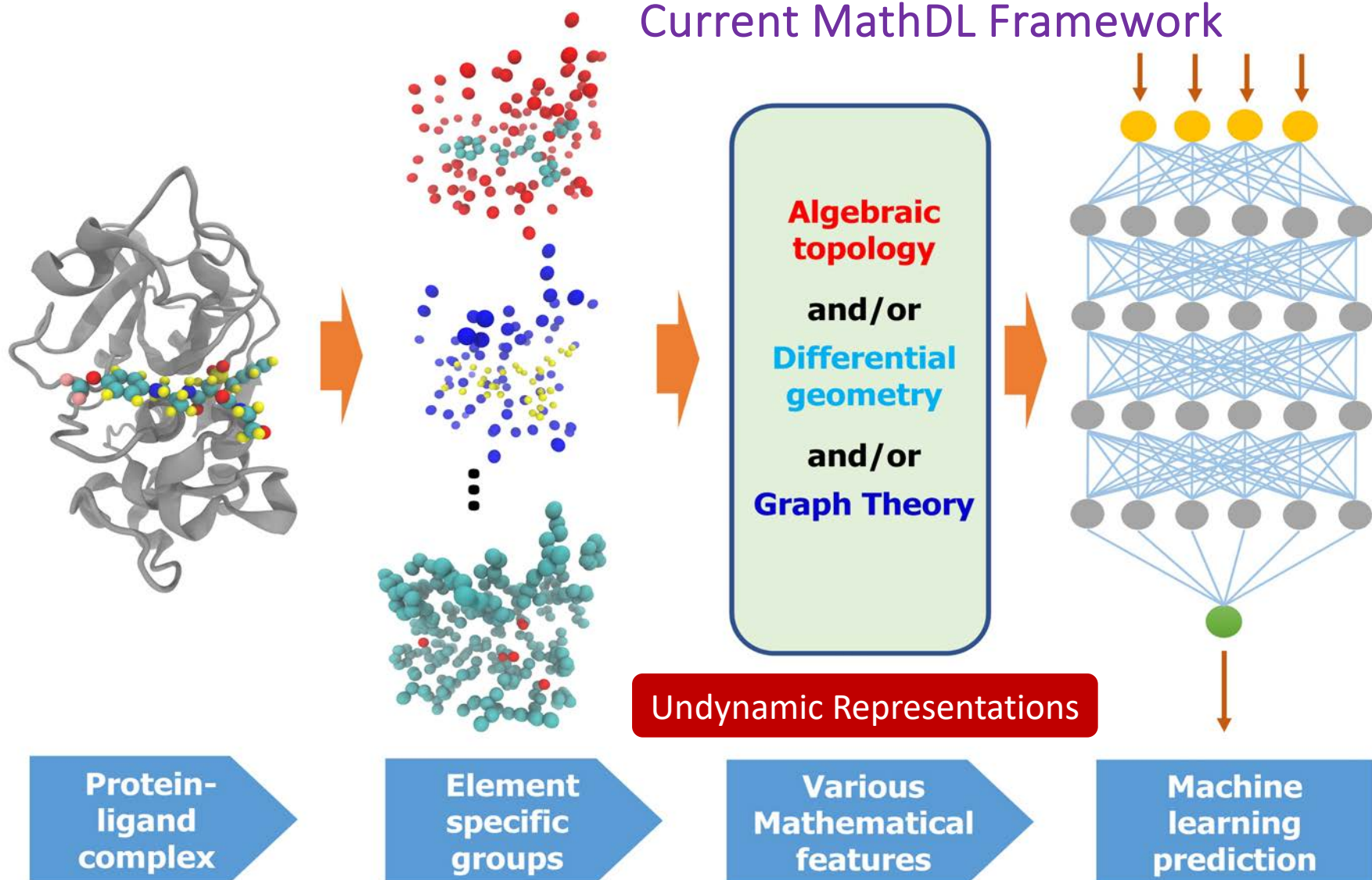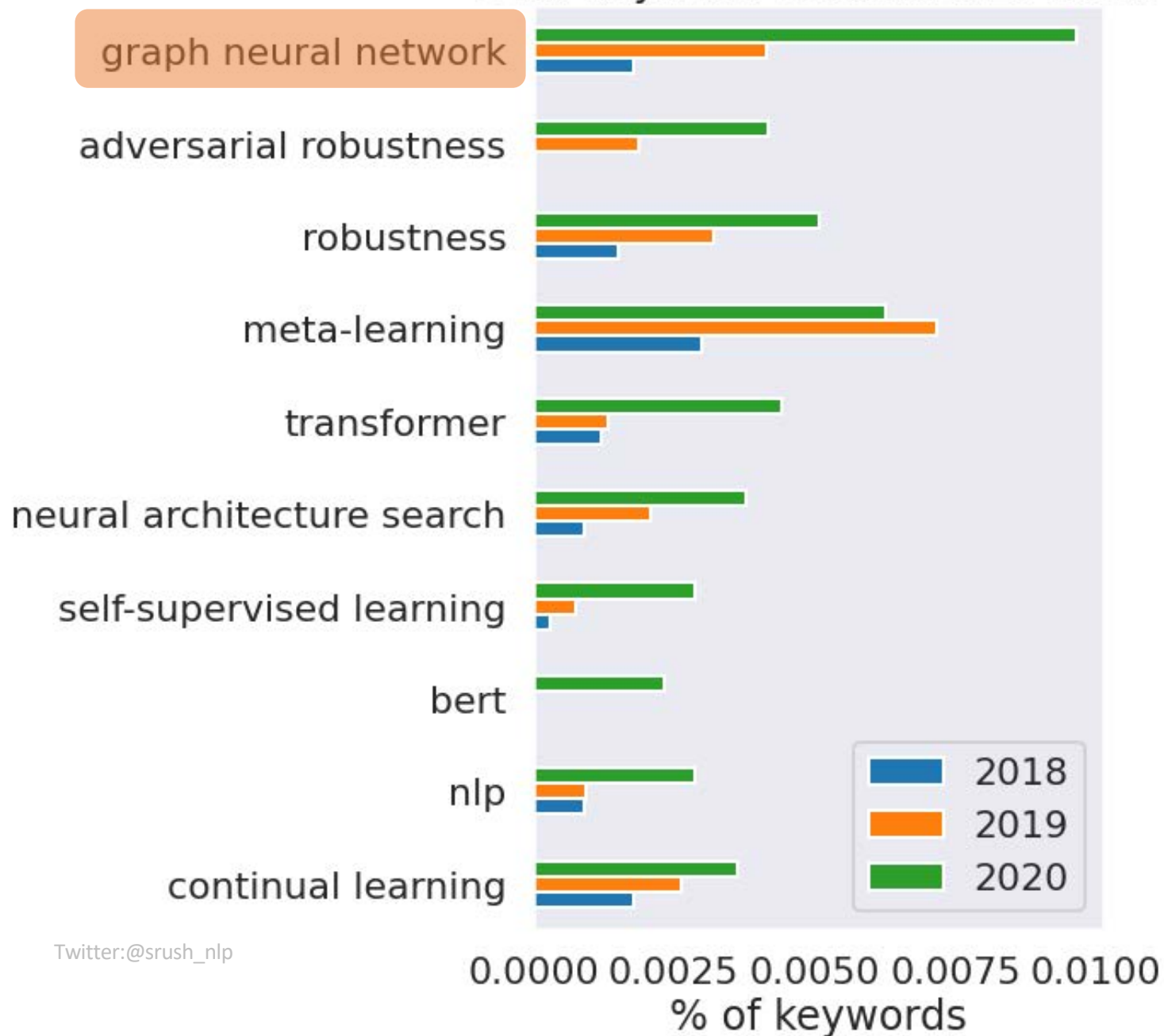Structure Based Scoring (Partials)
Free Energy Set 🥈🥉

(Nguyen et. al., JCAMD 2019)

# Moving Beyond MathDL



Current MathDL Framework

Algebraic topology
and/or
Differential geometry
and/or
Graph Theory

Undynamic Representations

**Protein-ligand complex** → **Element specific groups** → **Various Mathematical features** → **Machine learning prediction**
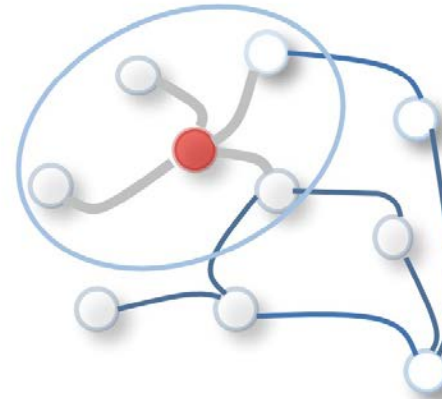
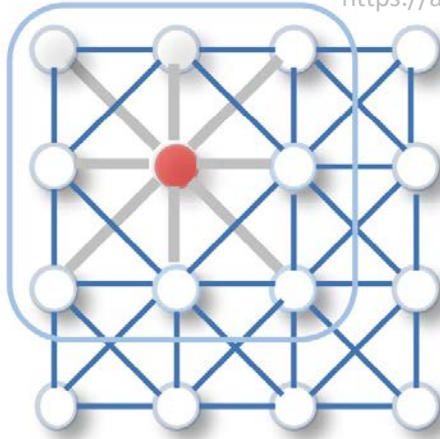ICLR Keyword Growth 2018-2020

Twitter:@srush_nlp

# Graph Neural Network (GNN)

Convolutional Layer

Graph Convolutional Layer

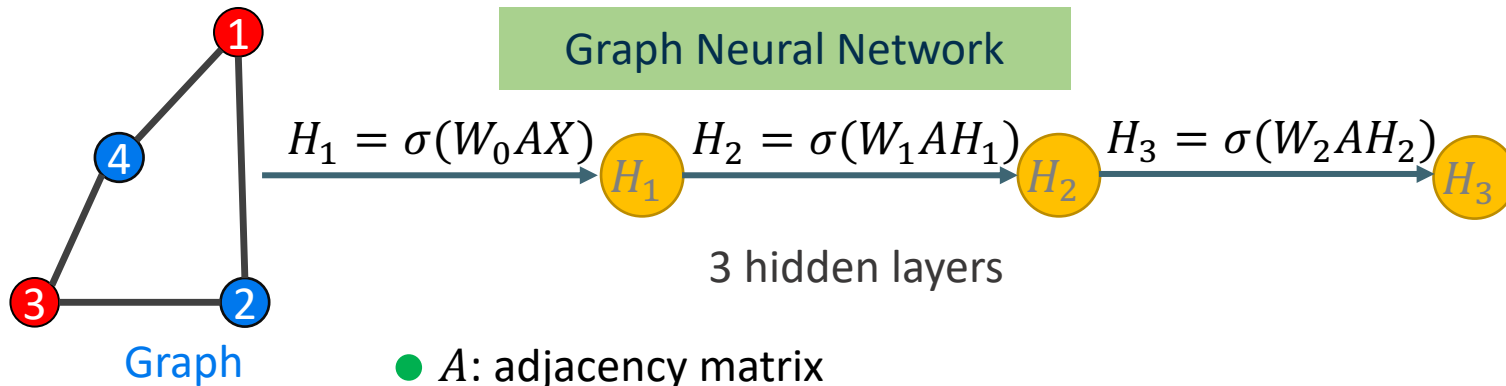https://arxiv.org/pdf/1901.00596.pdf

Deep Neural Network

$$z_1 = \sigma(W_0 x) \qquad z_2 = \sigma(W_1 z_1) \qquad z_3 = \sigma(W_2 z_2)$$

$x$    $z_1$    $z_2$    $z_3$

Node

3 hidden layers

Graph Neural Network

$$H_1 = \sigma(W_0 A X) \qquad H_2 = \sigma(W_1 A H_1) \qquad H_3 = \sigma(W_2 A H_2)$$

$H_1$    $H_2$    $H_3$

3 hidden layers

Graph

● $A$: adjacency matrix

# Standard Architecture of GNN



$\times K$

Read-Out

| Input | Graph Convolutional Network | Output Graph | Graph Representation | Fully Connected and Output |

# Standard Architecture of GNN

**● Aggregation:**

$$a_u^{(k)} = \text{AGGREGATE}^{(k)} \left( \left\{ h_v^{(k-1)}, v \in \mathcal{N}(u) \right\} \right)$$

Node State

Neighbor of node $u$

Example

$$a_u^{(k)} = \frac{1}{|\mathcal{N}(u)|} \sum_{v \in \mathcal{N}(u)} h_v^{(k-1)}$$

**● Combination/Updating:**

$$h_u^{(k)} = \text{COMBINE}^{(k)} \left( h_u^{(k-1)}, a_u^{(k)} \right)$$

Example

$$h_u^{(k)} = W_0^{(k)} h_u^{(k-1)} + W_1^{(k)} a_u^{(k)}$$
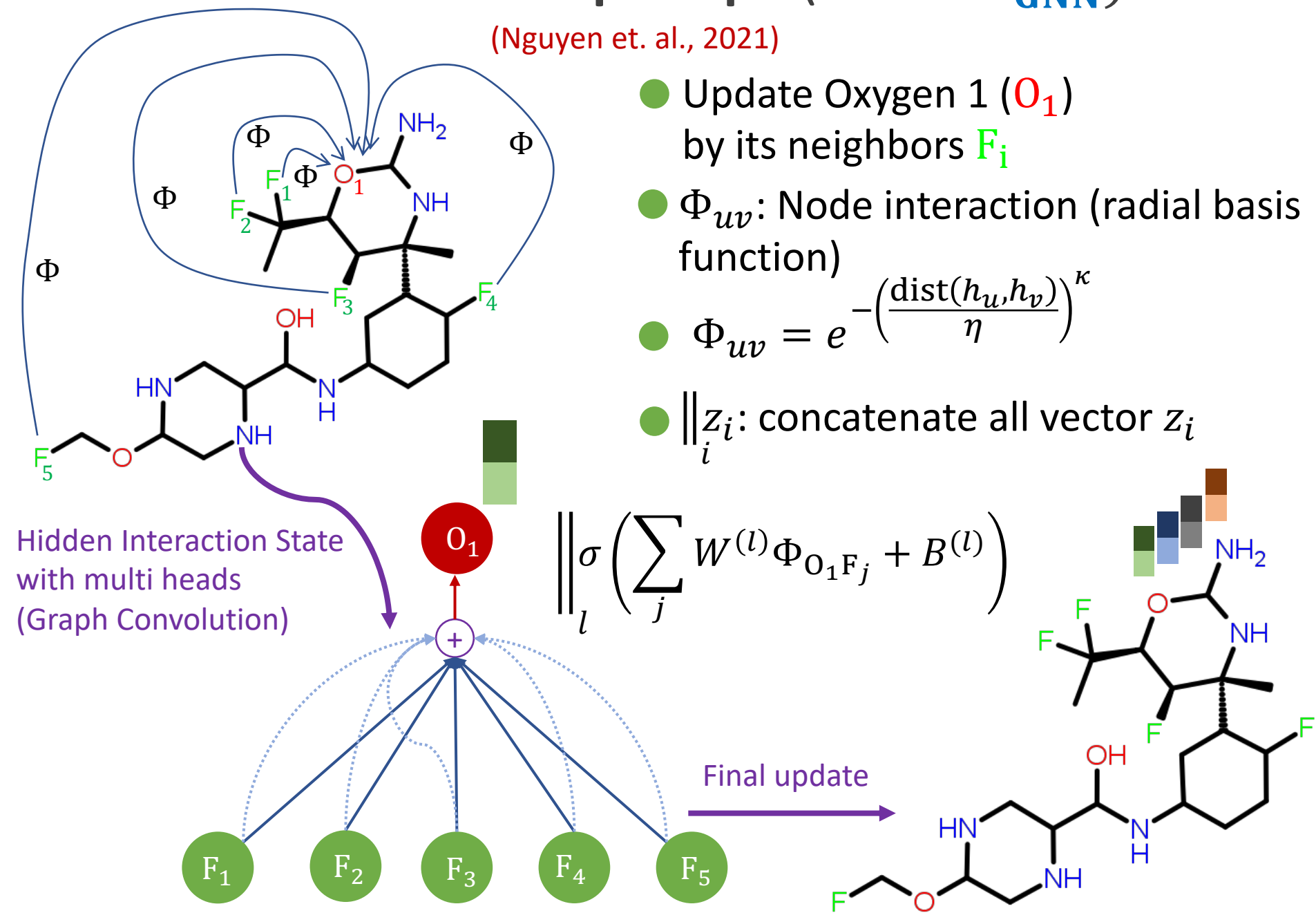
**● Read-out (Graph Invariant):**

$$h_G = \text{READOUT} \left( \left\{ h_u^{(K)}, u \in G \right\} \right)$$

Example

$$h_G = \text{MEAN} \left( \left\{ h_u^{(K)}, u \in G \right\} \right)$$

# Math-Based Deep Graph (MathDL$_{\text{GNN}}$)

(Nguyen et. al., 2021)



- Update Oxygen 1 ($O_1$) by its neighbors $F_i$

- $\Phi_{uv}$: Node interaction (radial basis function)

- $\Phi_{uv} = e^{-\left(\frac{\text{dist}(h_u, h_v)}{\eta}\right)^{\kappa}}$

- $\left\Vert_i z_i\right.$: concatenate all vector $z_i$

Hidden Interaction State with multi heads (Graph Convolution)

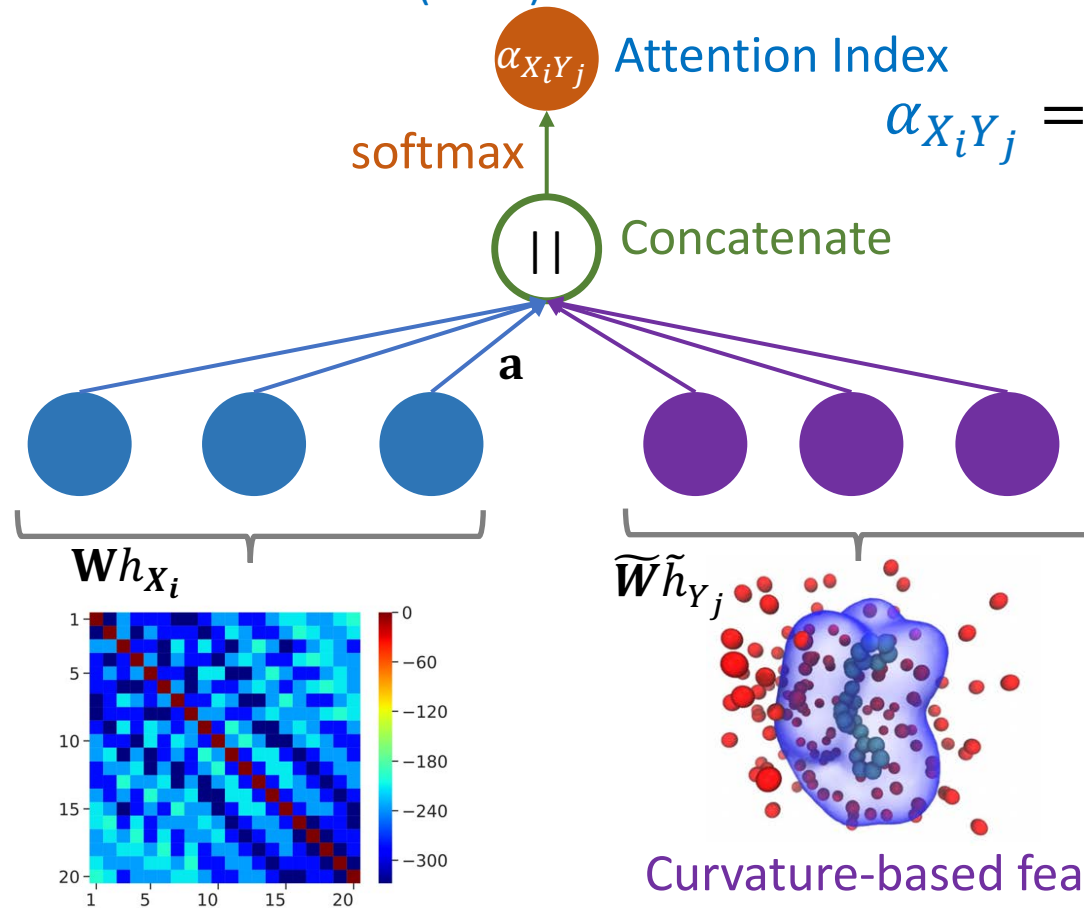$$\left\Vert_l \sigma\left(\sum_j W^{(l)}\Phi_{O_1 F_j} + B^{(l)}\right)\right.$$

Final update

# Math-Based Deep Graph (**MathDL$_{GNN}$**)

● **Attention mechanism** with another Math Features

Vaswani et. al. (2017)

$\alpha_{X_i Y_j}$ Attention Index

softmax

$||$ Concatenate

**a**

$$\alpha_{X_i Y_j} = \frac{\exp\left(\sigma(\mathbf{a}\left[\mathbf{W}h_{X_i}||\widetilde{\boldsymbol{W}}\tilde{h}_{Y_j}\right])\right)}{\sum_k \exp\left(\sigma(\mathbf{a}\left[\mathbf{W}h_{X_i}||\widetilde{\boldsymbol{W}}\tilde{h}_{Y_k}\right])\right)}$$

● New hidden state update

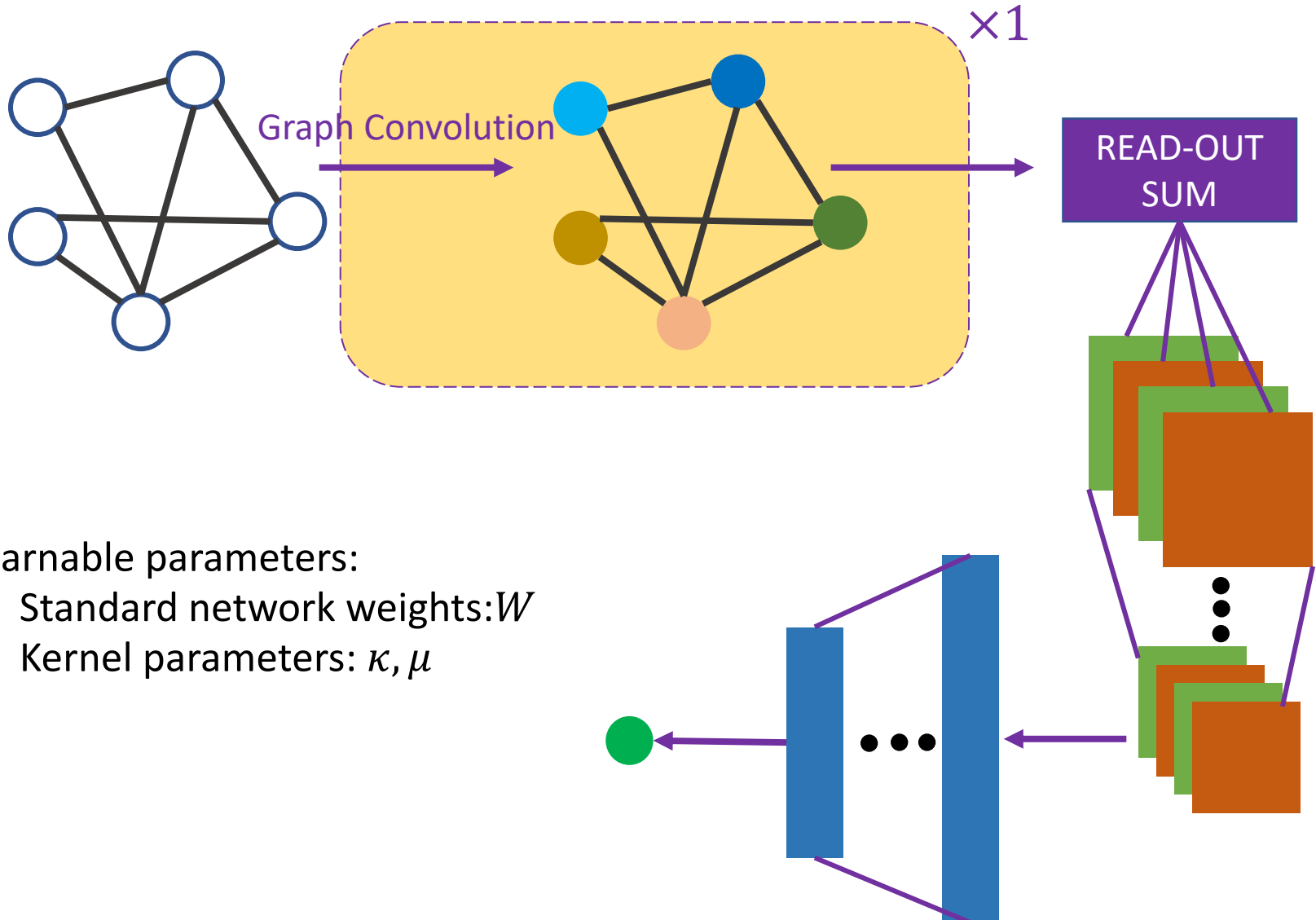$\mathbf{W}h_{X_i}$

$\widetilde{\boldsymbol{W}}\tilde{h}_{Y_j}$

$$h_{X_i}^{(new)} = \left\| \sigma\left( \sum_j \alpha_{X_i Y_j} W^{(k)} \Phi_{X_i Y_j} + B^{(k)} \right) \right.$$

Graph-based features

Curvature-based features
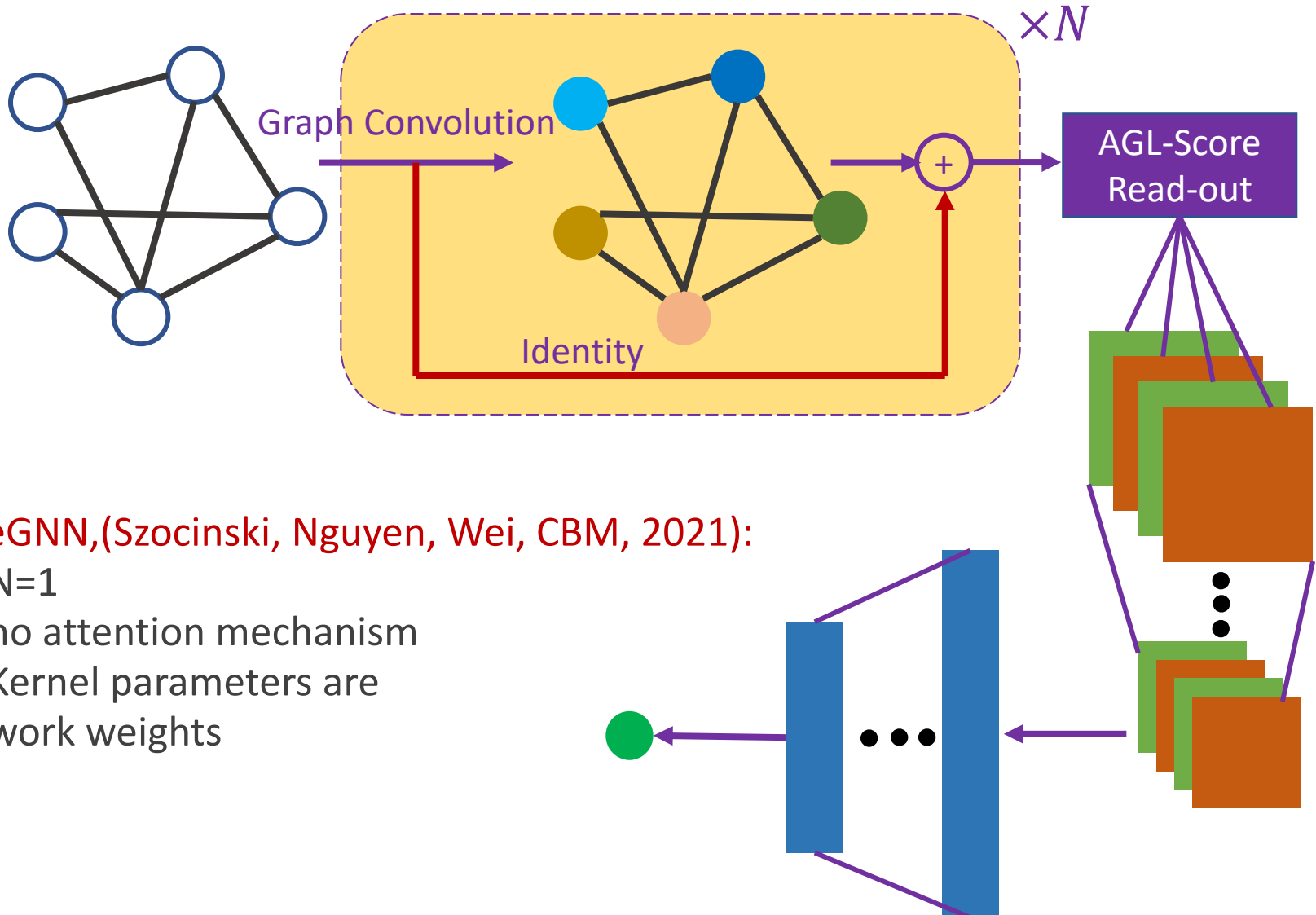(volume, area, electrostatics,
other physical interactions …)

# AweGNN

(Szocinski, Nguyen, Wei, CBM, 2021)

$\times 1$

Graph Convolution

READ-OUT SUM

- Learnable parameters:
  - Standard network weights: $W$
  - Kernel parameters: $\kappa, \mu$

# Math-Based Deep Graph (MathDL$_{GNN}$)

(Nguyen et. al., 2021)

$\times N$

Graph Convolution

Identity

+

AGL-Score Read-out
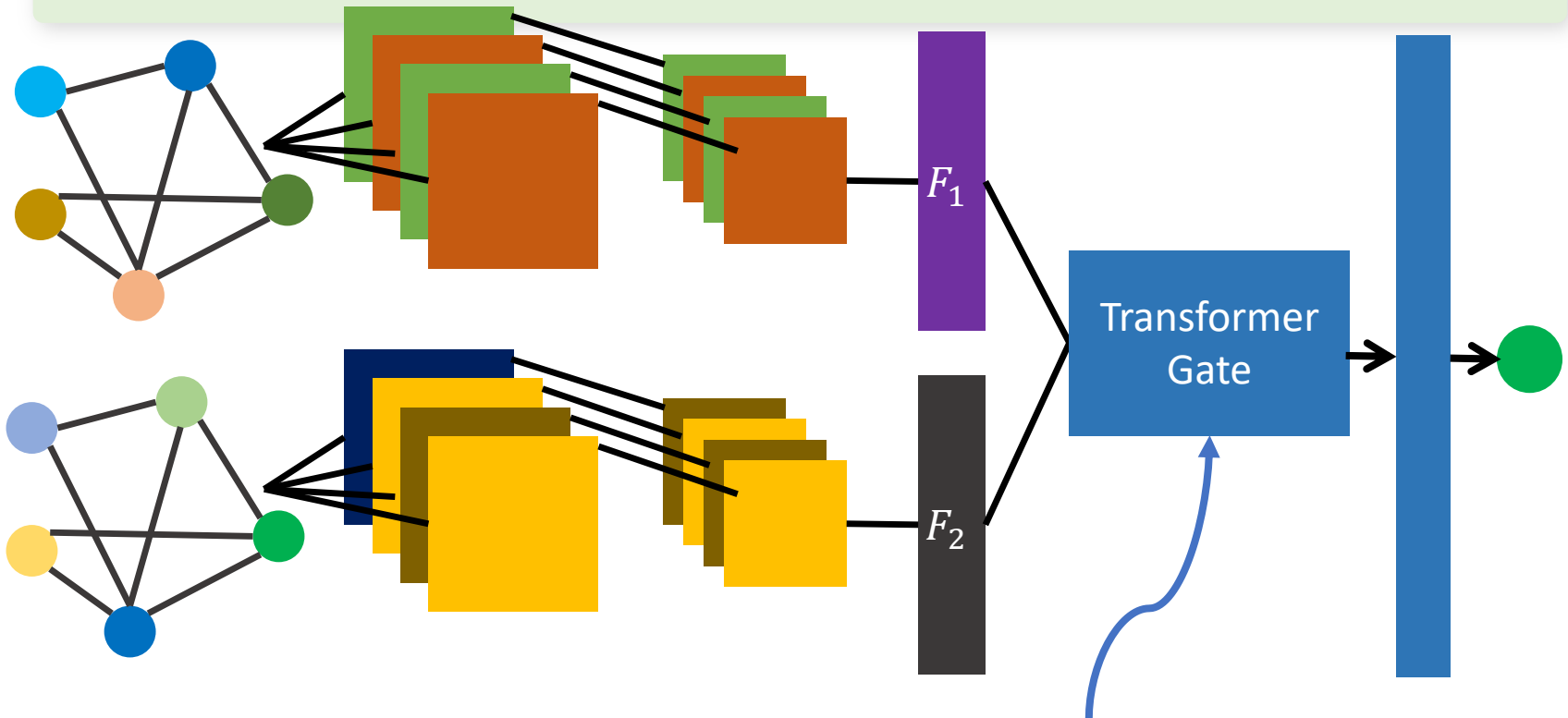
● AweGNN,(Szocinski, Nguyen, Wei, CBM, 2021):
- N=1
- no attention mechanism
- Kernel parameters are network weights

# Math-Based Deep Graph (**MathDL**$_{\text{GNN}}$)

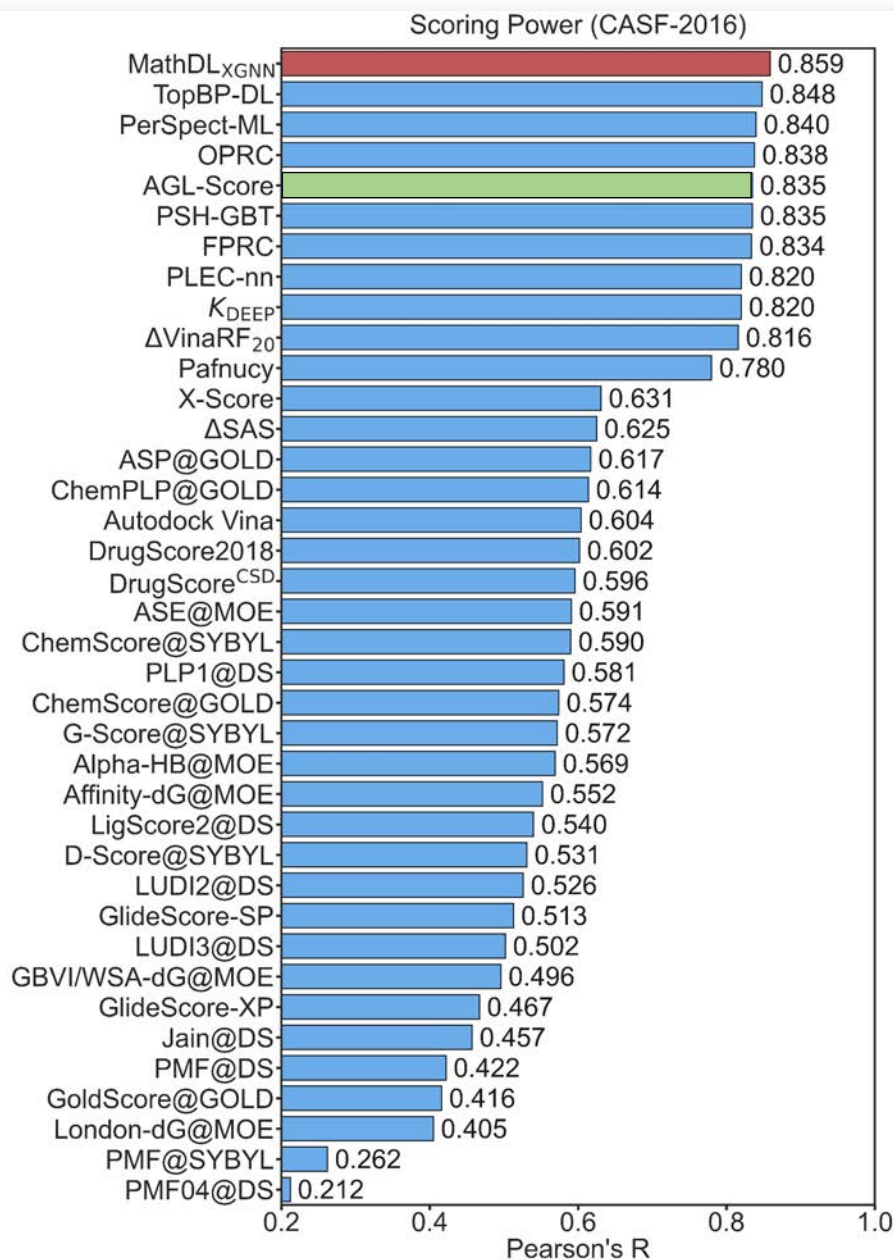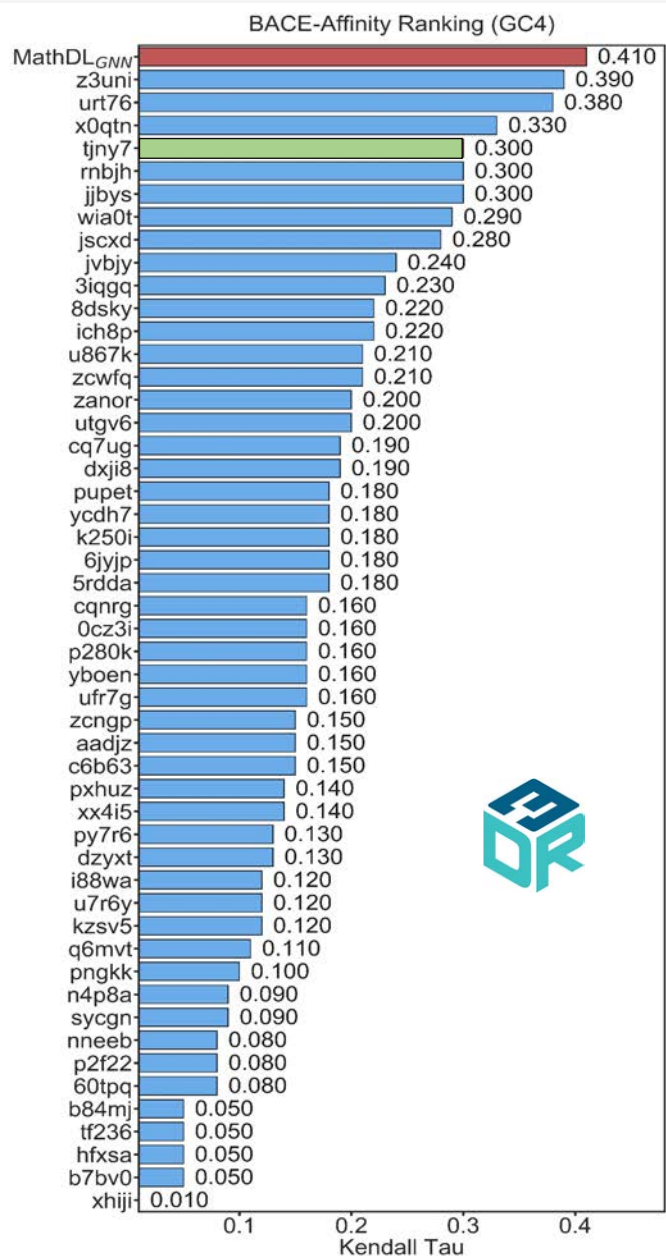● **Transformer gate:** combining different Math-based Graphs



$F_1$

$F_2$

Transformer Gate

(Nguyen et. al., 2021)

$$z = \sigma\left(\widetilde{W}_1 F_1 + \widetilde{W}_2 F_2 + b\right)$$

$$F_{\text{combined}} = z \odot W_1 F_1 + (1 - z) \odot W_2 F_2$$

# Performance of MathDL in Scoring Power

(Nguyen et. al., 2021)

Thank you!