

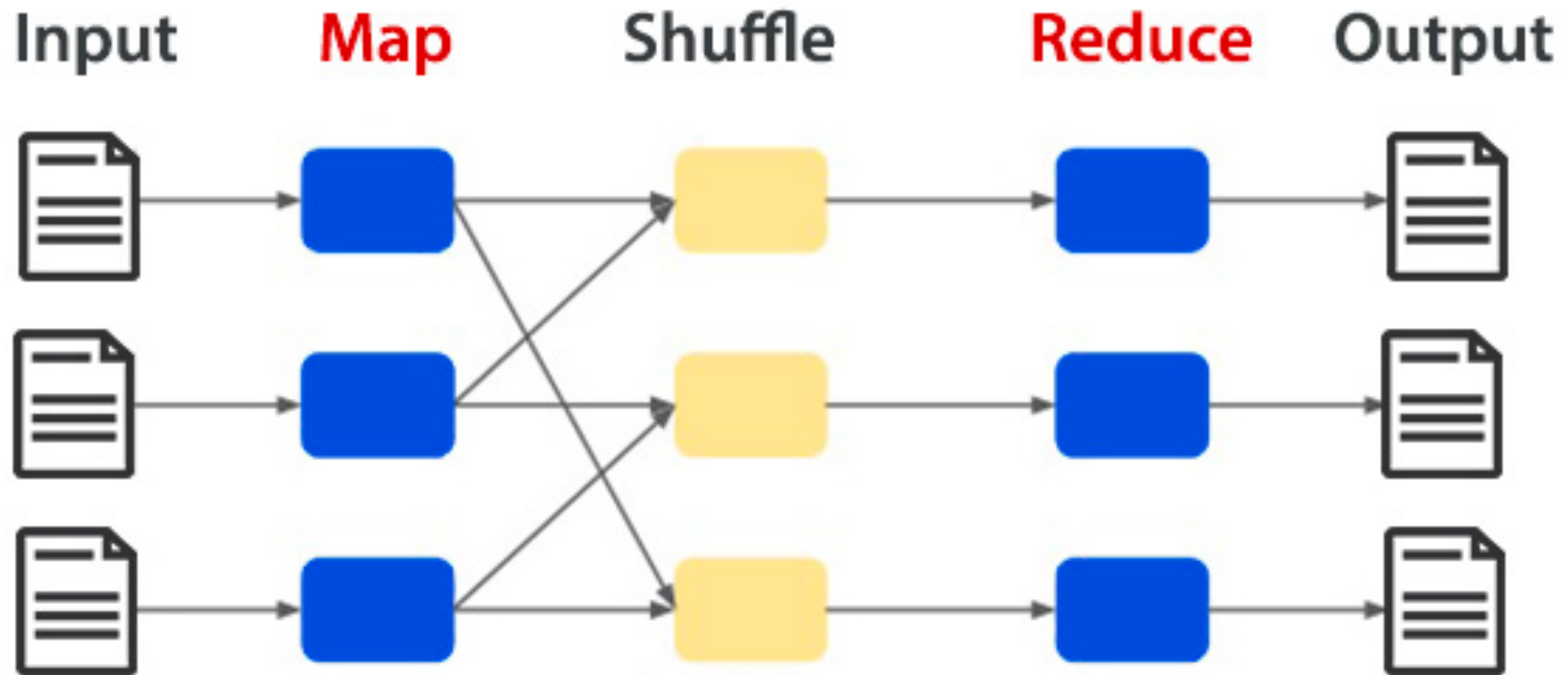
STEFANIE MUROYA LEI Y JUANPABLO HEREDIA PARILLO

TOWARDS QUANTUM FEATURE SELECTION APPLIED TO MASSIVE DATASETS

Desarrollo de un método que pueda sere escalable y útil para datasets con las dimensiones usadas en Big Data.

MARCO TEÓRICO

BIG DATA: MAPREDUCE



@edupristine.com

TABU SEARCH

- Método meta-heurístico de búsqueda.
- Métodos de búsqueda local son usados para optimización matemática.

```
BEGIN  
   $t \leftarrow 0$ ;  
  INITIALIZE TABU SEARCH;  
  WHILE ( $t < t_{max}$ ) DO:  
     $t \leftarrow t + 1$ ;  
    SEARCH NEIGHBORHOOD;  
    EVALUATE CANDIDATE SOLUTIONS;  
    UPDATE TABU LIST;  
  END  
END
```

TEORÍA DE LA INFORMACIÓN

- Teoría matemática que define los límites y posibilidades de la comunicación.
- Se basa de teoría de probabilidades y en estadística.
- Una medida clave en teoría de la información es la **entropía**.

$$H(X) = - \sum_x p(x) \log_2 p(x)$$

ENTROPÍA

$$H(X, Y) = - \sum_{x,y} p(x, y) \log_2 p(x, y)$$

ENTROPÍA CONJUNTA

$$H(X|Y) = - \sum_{x,y} p(x,y) \log_2 p(x|y)$$

ENTROPÍA CONDICIONAL

INFORMACIÓN MUTUA (IM)

Mide la cantidad de información que se puede obtener sobre una variable aleatoria observando a otra.

$$I(X; Y) = H(X) - H(X|Y)$$

INFORMACIÓN MUTUA CONDICIONAL (IMC)

Es el valor esperado de la información mutua entre 2 variables aleatorias dada una tercera.

$$I(X; Y|Z) = H(X|Z) - H(X|Y, Z)$$

MECÁNICA CUÁNTICA

$$|0\rangle = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$|1\rangle = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$|0\rangle = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$|1\rangle = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$|q_1\rangle = \alpha_1|0\rangle + \beta_1|1\rangle$$

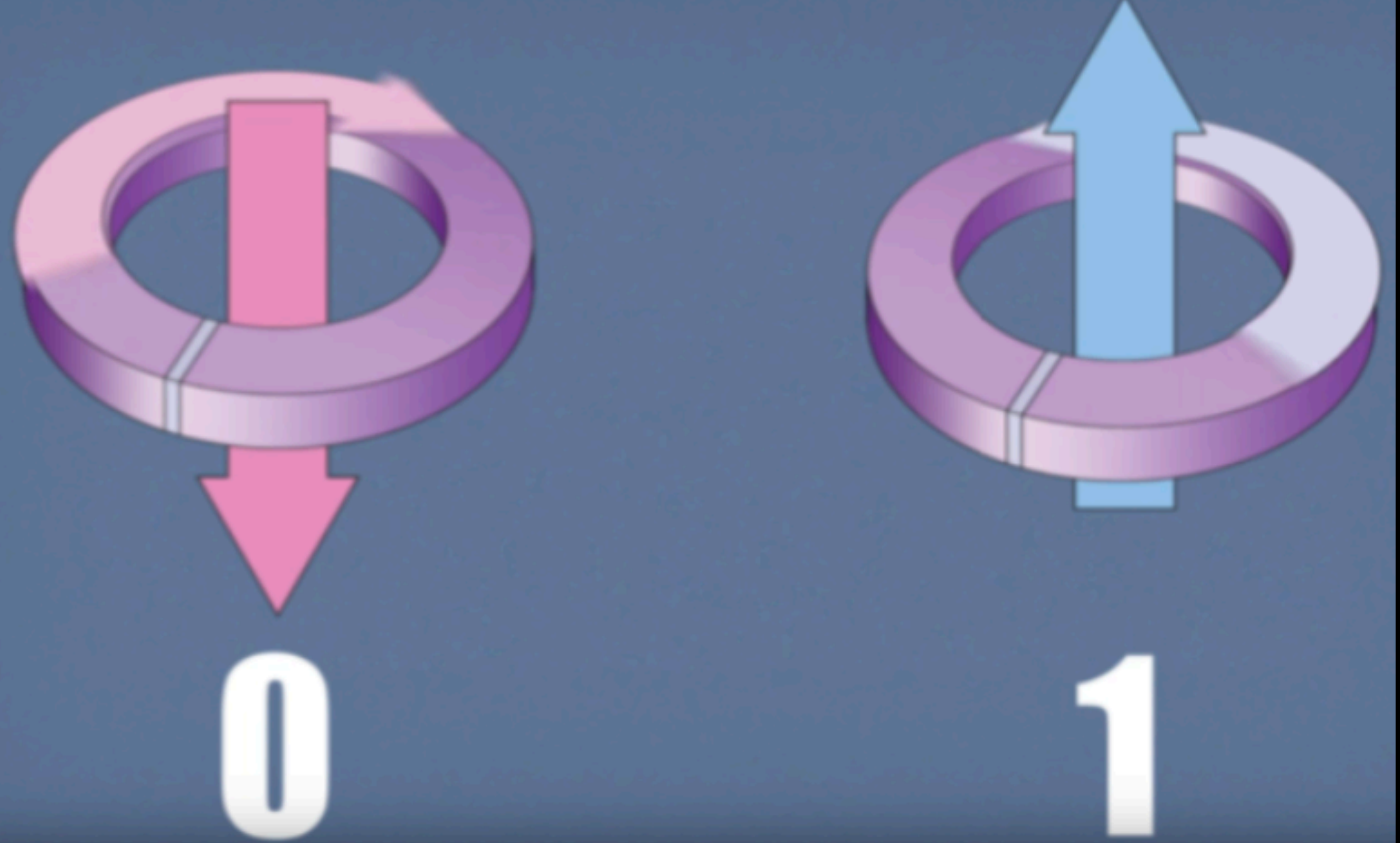
$$|q_2\rangle = \alpha_2|0\rangle + \beta_2|1\rangle$$

$$|q_1q_2\rangle = \alpha_1\alpha_2|00\rangle + \alpha_1\beta_2|01\rangle + \beta_1\alpha_2|10\rangle + \beta_1\beta_2|11\rangle$$

“UN SISTEMA CUÁNTICO ES UN PRODUCTO DE ESTADOS SI ES QUE EXISTE UNA FORMA DE ESCRIBIRLO COMO UN PRODUCTO DE TENSORES DE 1-QUBIT, EN CASO CONTRARIO, ES UN SISTEMA ENTRELAZADO.”

ENTRELAZAMIENTO DE QUBITS

QUANTUM ANNEALING



TOPOLOGICAL

GATE MODEL

MEASUREMENT BASED

**UNIVERSAL QUANTUM
COMPUTING**

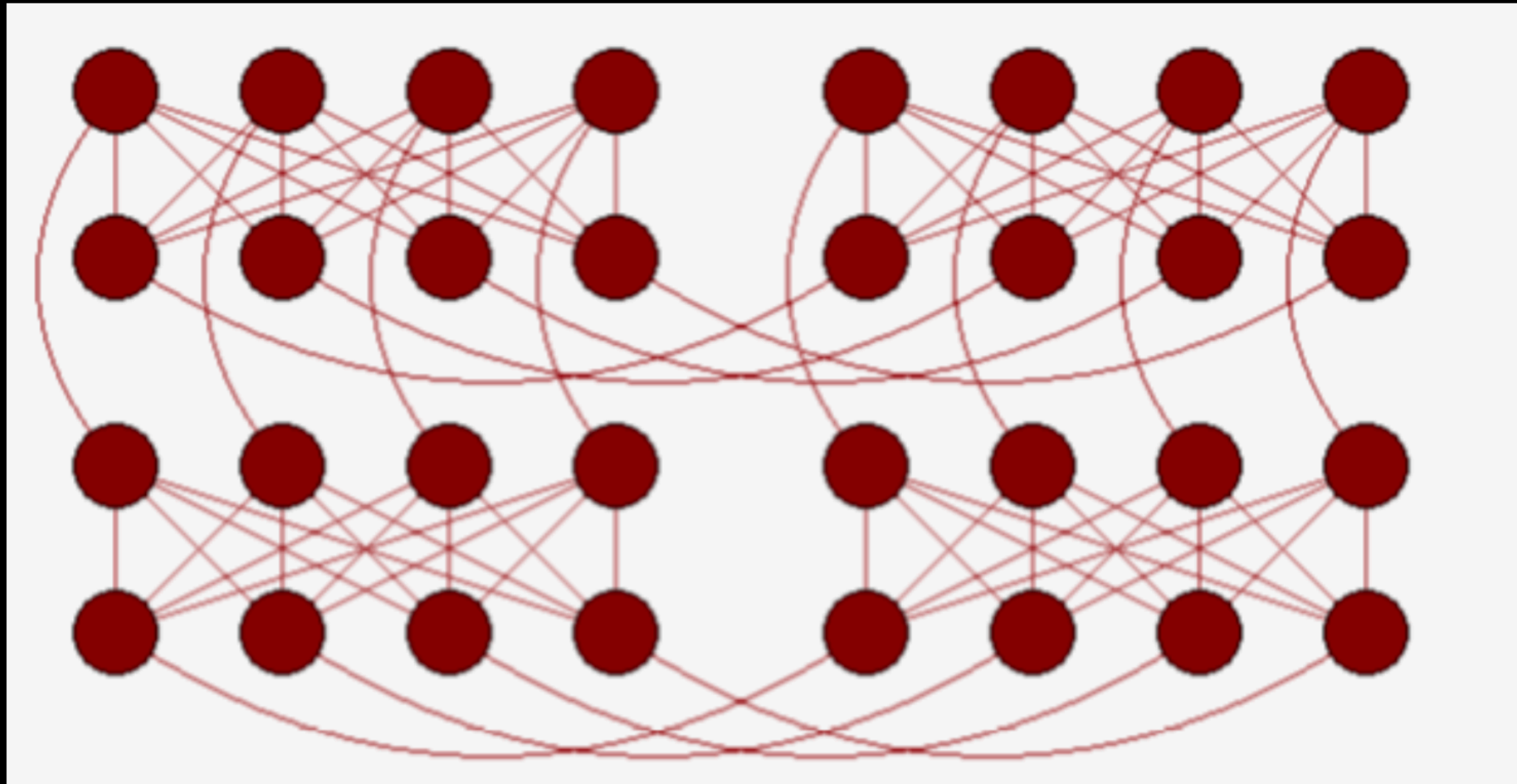
ADIABATIC

DWAVE

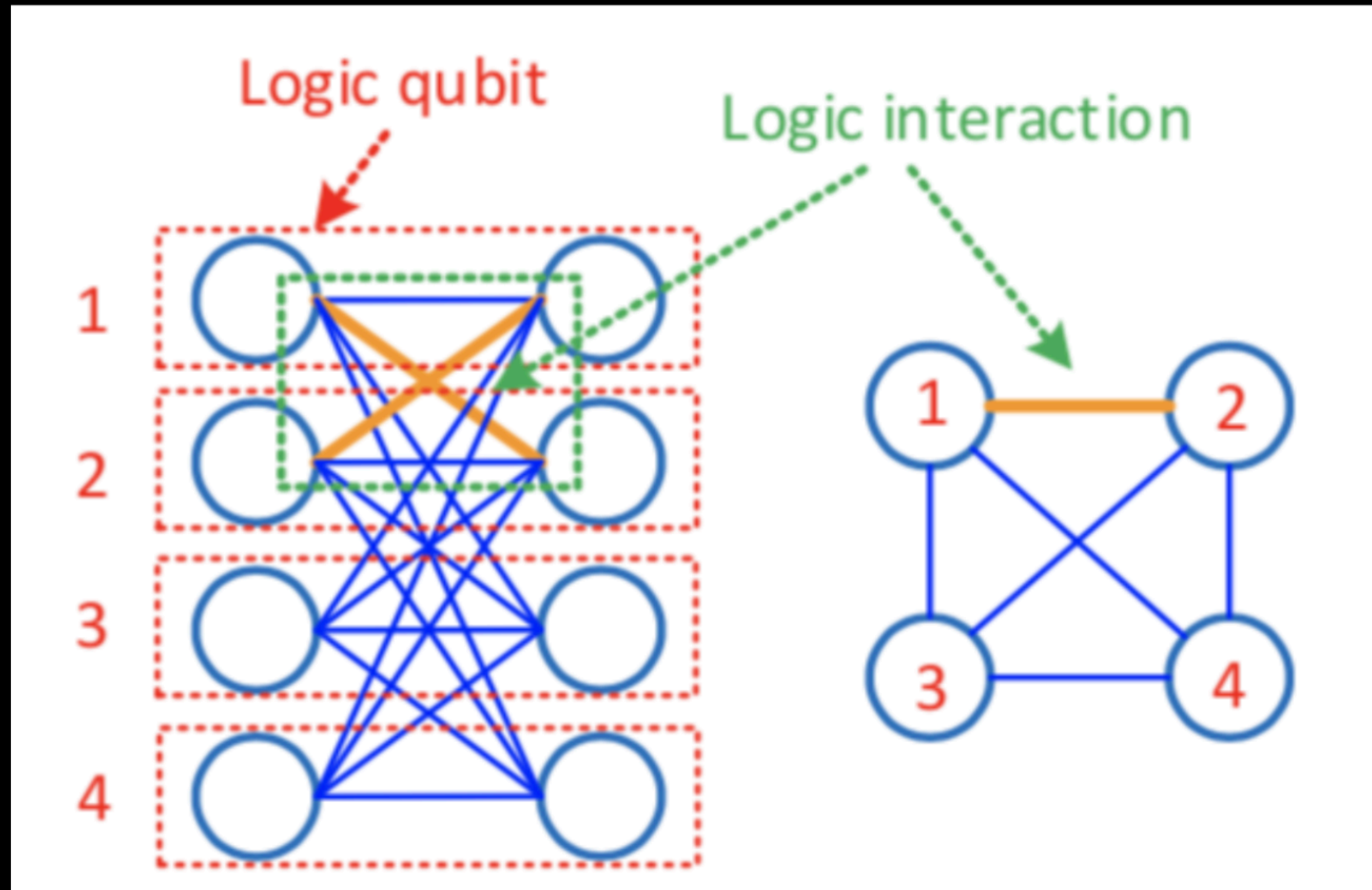
QUANTUM ANNEALING

ARQUITECTURA Y DISPOSICIÓN DE QUBITS

$$I = \sum_i^N h_i s_i + \sum_{i=1}^N \sum_{j=i+1}^N J_{i,j} s_i s_j$$

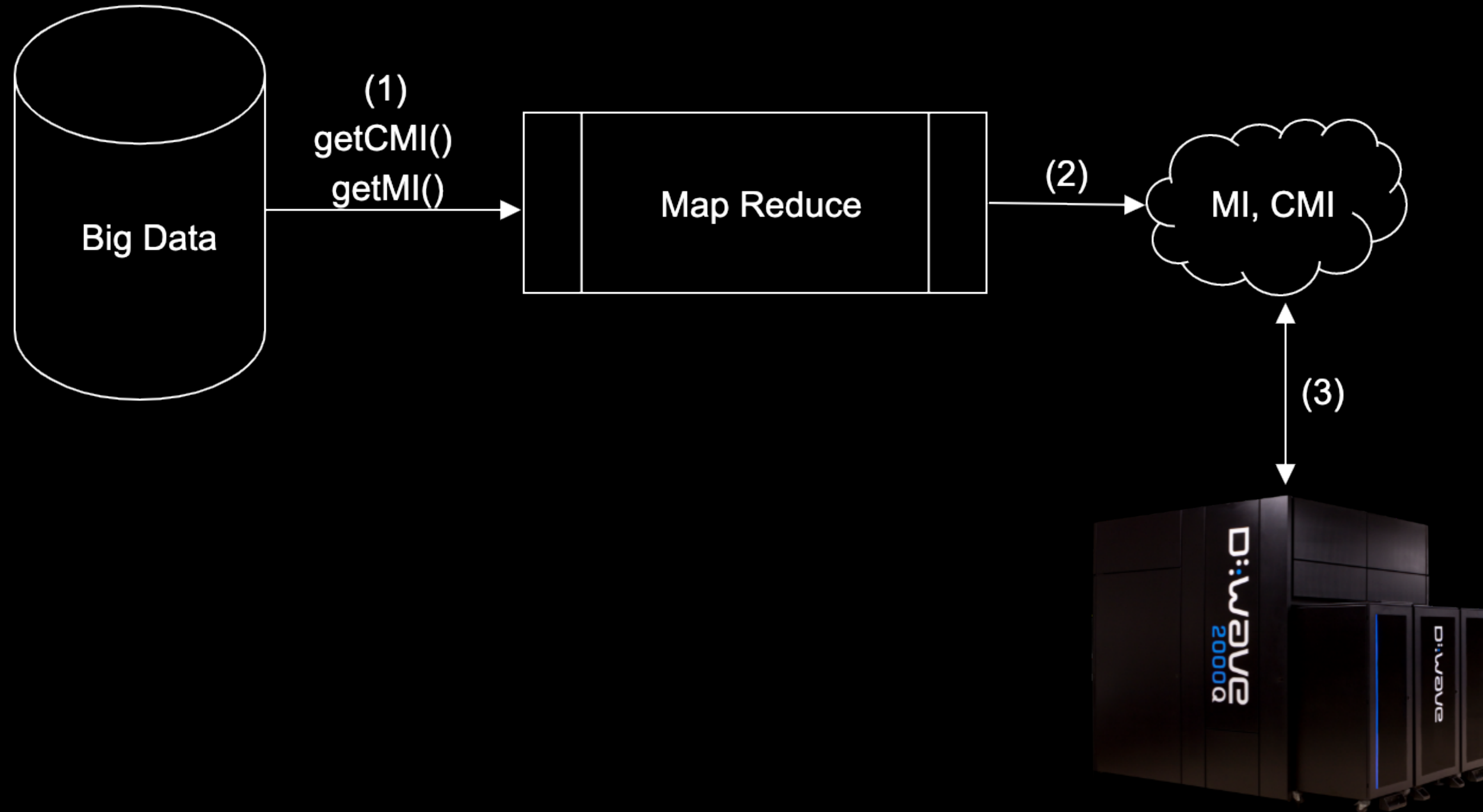


QUBITS LOGICOS VS. QUBITS FÍSICOS



PROPUESTA

ESTRUCTURA PRINCIPAL



JUSTIFICACIÓN

$$\mathcal{O}(IM \ \& \ IMC) + \mathcal{O}(k_feat \ . \ _selec) = \mathcal{O}(F^2) + \mathcal{O}(F!)$$

$$QUBO = \sum_i^n IM(X_i, Y)x_i + \sum_i \sum_j ICM(X_j; Y | X_i)x_i x_j$$

PENALIZANDO RESPUESTAS INVÁLIDAS

$$\alpha \sum_i^n (x_i - k)^2$$

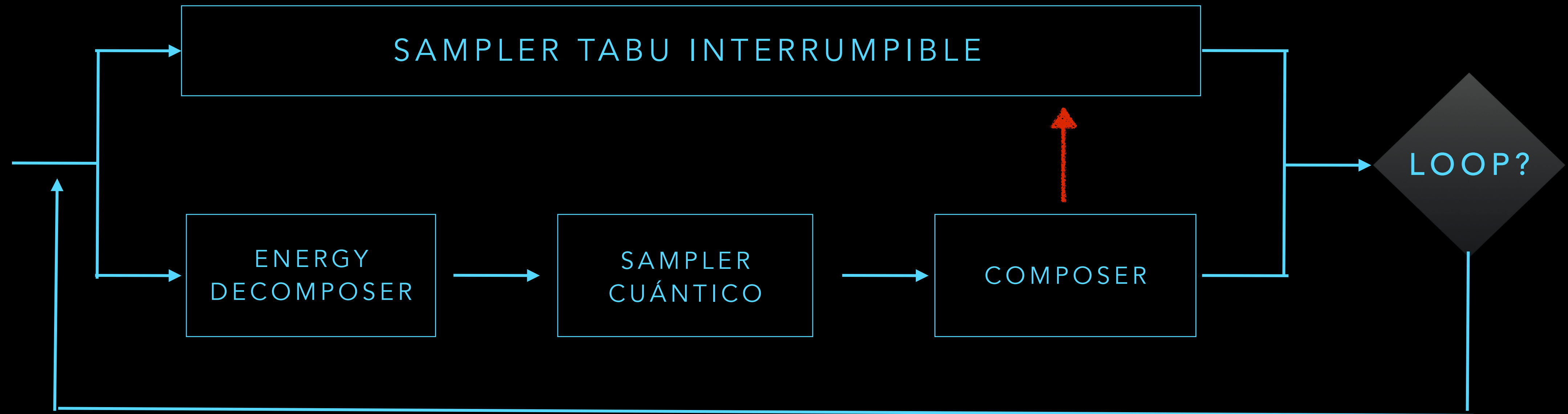
BIG DATA: HYBRID QUANTUM SOLVER

RAMAS

INTERRUPTABLE_TABU_SAMPLER()

DECOMPOSER | SAMPLER | COMPOSER

WORKFLOW PRINCIPAL



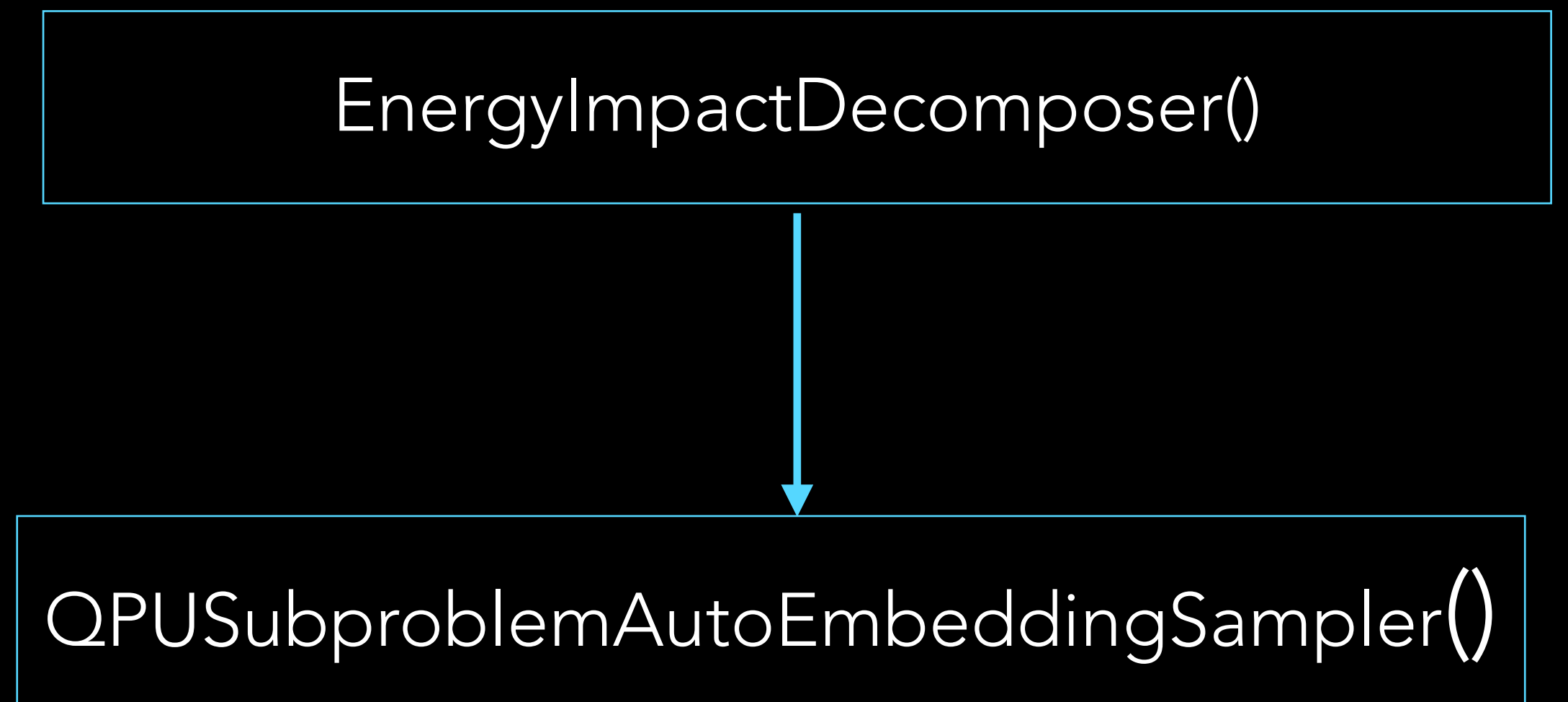
ENERGY IMPACT DECOMPOSER

- Crea subproblemas de un determinado tamaño (batches).

```
EnergyImpactDecomposer()
```

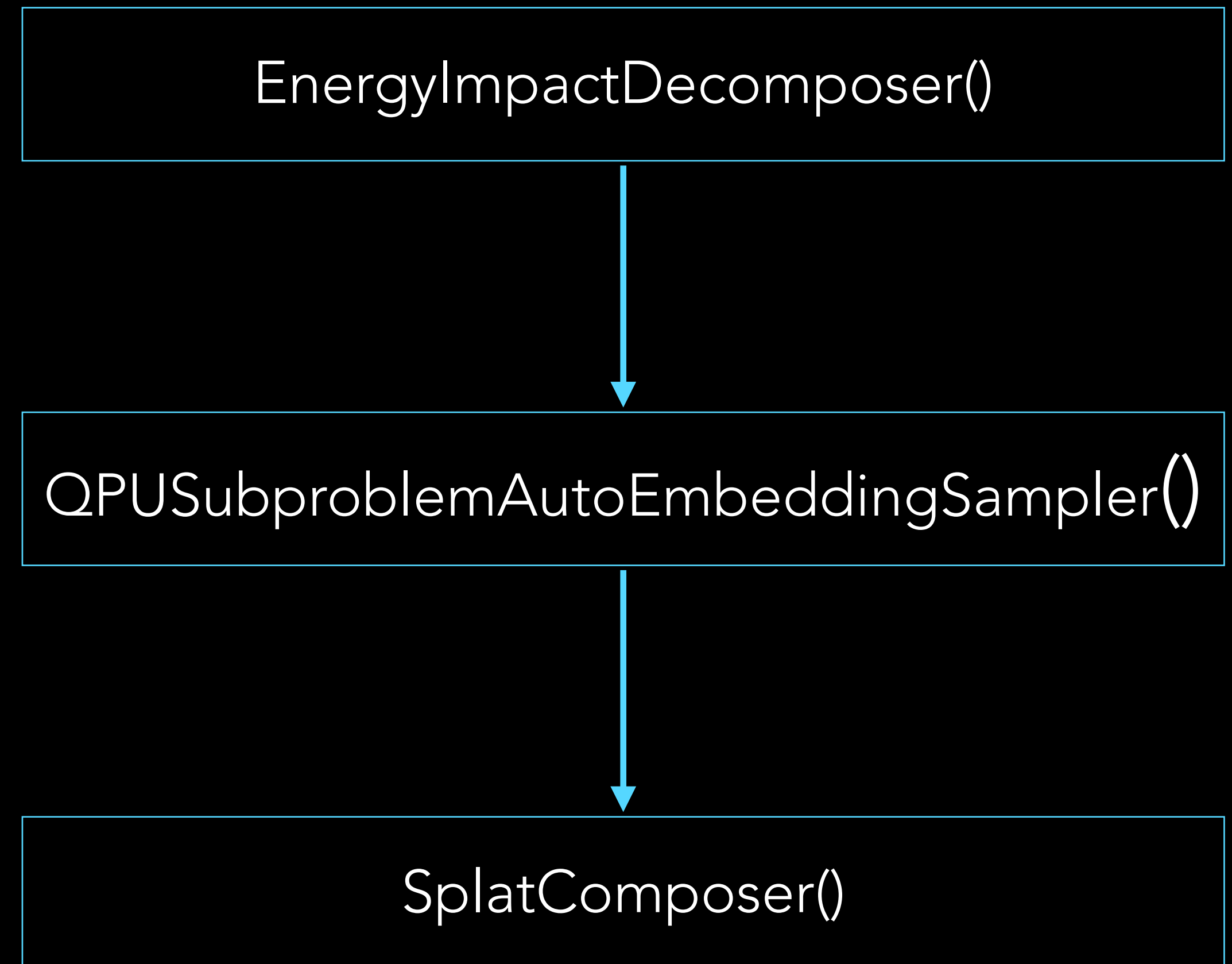
QUANTUM SAMPLER

- Ejecuta un sub-problema en la computadora cuántica



COMPOSER

- Reemplaza y actualiza las soluciones tabu.



EXPERIMENTOS

HARDWARE

CLÁSICO

- 2.3 GHz Dual-Core Intel Core i5
- 8 GB LPDDR3
- 4 núcleos físicos
- 8 núcleos lógicos

CUÁNTICO (DW_2000Q_6)

- 2041 qubits
- 1-10s tiempo de espera.
- **Total post-procesamiento:** $560\ \mu s$
- **Tiempo de iteración (anneal):** $20\ \mu s$
- **Tiempo programación qpu:** $10719\ \mu s$
- **Tiempo de lectura de qubits (por iteración):** $198\ \mu s$

SOFTWARE

CLÁSICO

- Python 3.7.6

CUÁNTICO

- **Dwave-system** 0.9.6
- **Dimod** 0.9.4
- **dwave-hybrid** 0.5.0

DATASET 1: TITANIC

¿Cuáles son los mejores atributos para predecir si un pasajero sobrevivió?

1045 registros

14 atributos

batch_size = 2

$\alpha : 10$

Qubits por Variable: 5

RESULTADO PARA K=3

0.578 segundos en QPU

¿Género?

¿ Mr.?

¿ Famoso?

$$\sum IM + \sum IMC = 1.46$$

RESPUESTA

Género

Mr.

Miss.

$$\sum IM + \sum IMC = 2.48$$

DATASET 2: HOUSING

¿Cuáles son los mejores atributos para predecir el precio de una casa?

1460 registros

batch_size: 10

80 atributos

convergencia: 3

PARA $K = 20$, $\alpha = 100$

~6 SEG

1. MSSubClass	8. Exterior1st	15.2ndFlrSF
2. LotFrontage	9. Exterior2nd	16.GrLivArea
3. LotArea	10.BsmtFinType1	17.GarageYrBlt
4. Utilities	11.BsmtFinSF1	18.GarageArea
5. Neighborhood	12.BsmtUnfSF	19.OpenPorchSF
6. YearBuilt	13.TotalBsmtSF	20.MoSold
7. YearRemodAdd	14.1stFlrSF	21.YrSold

PARA $K = 20, \alpha = 100000$

~6 SEG

- | | | |
|-----------------|-----------------|----------------|
| 1. MSSubClass | 8. Exterior1st | 15.2ndFlrSF |
| 2. LotFrontage | 9. Exterior2nd | 16.GrLivArea |
| 3. LotArea | 10.BsmtFinType1 | 17.GarageYrBlt |
| 4. Utilities | 11.BsmtFinSF1 | 18.GarageArea |
| 5. Neighborhood | 12.BsmtUnfSF | 19.OpenPorchSF |
| 6. YearBuilt | 13.TotalBsmtSF | 20.YrSold |
| 7. YearRemodAdd | 14.1stFlrSF | |

CONCLUSIONES

- Se desarrolló un método escalable y fácil de implementar.
- Recalentamiento cuántico tiene un futuro prometedor.
- Se vieron limitaciones de hardware producto de una tecnología que es relativamente nueva en el ámbito empírico.
- Hemos demostrado supremacía cuántica.

TRABAJOS FUTUROS

- Mejorar la técnica para espectros poco amplios.
- Desarrollar componentes de workflow que mejoren el desempeño.

Bibliografía

- [HTTPS://CLOUD.DWAVESYS.COM/LEAP/](https://cloud.dwavesys.com/leap/)
- [HTTPS://DOCS.OCEAN.DWAVESYS.COM/EN/STABLE/DOCS_HYBRID/INTRO/USING.HTML](https://docs.ocean.dwavesys.com/en/stable/docs_hybrid/intro/using.html)
- [HTTPS://TOWARDSDATASCIENCE.COM/FEATURE-SELECTION-TECHNIQUES-IN-MACHINE-LEARNING-WITH-PYTHON-F24E7DA3F36E#:~:TEXT=FEATURE%20SELECTION%20IS%20THE%20PROCESS,LEARN%20BASED%20ON%20IRRELEVANT%20FEATURES](https://towardsdatascience.com/feature-selection-techniques-in-machine-learning-with-python-f24e7da3f36e#:~:text=FEATURE%20SELECTION%20IS%20THE%20PROCESS,LEARN%20BASED%20ON%20IRRELEVANT%20FEATURES).
- BROWN, G., POCOCK, A., ZHAO, M. J., & LUJÁN, M. (2012). "CONDITIONAL LIKELIHOOD MAXIMISATION: A UNIFYING FRAMEWORK FOR INFORMATION THEORETIC FEATURE SELECTION." THE JOURNAL OF MACHINE LEARNING RESEARCH, 13(1), 27-66.
- X. V. NGUYEN, J. CHAN, S. ROMANO, AND J. BAILEY, "EFFECTIVE GLOBAL APPROACHES FOR MUTUAL INFORMATION BASED FEATURE SELECTION". A. MONTANARO, "QUANTUM ALGORITHMS: AN OVERVIEW," NPJ QUANTUM INFORMATION, VOL. 2, P. 15023, JAN 2016.
- A. LEONARD SUSSKIND, QUANTUM MECHANICS.BASIC BOOKS, 2014.