# CHALLENGE

## Questions:

**Q1**.

Given the unique requirements and constraints of the challenge, we'll need to craft data structures that balance efficiency and compactness. The objective is to minimize storage while maintaining the fidelity of information contained in each image type. Let's dive into the specifics for both image types.

### Microscope Images

Data Structure: Since microscope images depict the parasites as blobs of arbitrary shapes in a binary color scheme (black for the parasite, white for the surrounding), a run-length encoding (RLE) scheme is suitable. This approach compresses sequences of data in which the same data value occurs in many consecutive data elements.

**Rationale:** RLE is efficient for this case because it can significantly reduce the size of large areas of uniform color. Given that each parasite occupies at least 25% of the image area, and considering the binary color scheme, RLE can efficiently compress these large sections of black or white.

**Storage Estimate:** For a 100,000x100,000 pixel image, there are 10 billion pixels. In the worst case, alternating black and white pixels would result in the least compression, but this is highly unlikely given the described scenario. More realistically, large areas of uniform color would allow for substantial compression. Assuming an average run length of 100 pixels (a conservative estimate given the described scenarios), and each run being described by a 32-bit integer for the run length and a 1-bit value for the color (black or white), the storage per run would be 4.125 bytes. Thus, for an image, the worst-case scenario might be closer to $10^{10}/100 \times 4.125$ bytes.

### Dye Sensor Images

Data Structure: A sparse matrix representation, specifically a Coordinate List (COO). This format stores only the locations and values of non-zero elements, which in this case would be the pixels where the dye is lit up.

**Rationale**: Given that the dye is expected to light up specific parts of the parasite and there's some leakage outside, most of the image would still be black (no dye), making it sparse. The

COO format is ideal for storing sparse data where the non-zero elements are scattered, as it doesn't waste space on the zero elements.

**Storage Estimate**: Assuming fewer than 0.1% of the parasites have cancer, and thus a similarly low percentage of the image area would actually contain lit dye (non-zero elements), the storage required would be much less than for a full bitmap representation. Each non-zero pixel would be represented by its two-dimensional coordinate (each a 32-bit integer) and a 1-bit value for the lit dye, totaling 8.125 bytes per non-zero pixel. For an image with 0.1% of its pixels lit, that would be $10^{10} \times 0.001 \times 8.125$ bytes.

**Worst-case Storage Size**

- Microscope Images: Worst-case (with minimal compression): $10^{10}/100 \times 4.125$ bytes.
- Dye Sensor Images: Assuming 0.1% of pixels are non-zero (high estimate for cancerous cases): $10^{10} \times 0.001 \times 8.125$ bytes.

Let's calculate these estimates to provide a clearer answer:
The worst-case storage size estimates for each image representation are as follows:
- Microscope Images: Approximately 412.5 million bytes (or roughly 393 MB) per image.
- Dye Sensor Images: Approximately 81.25 million bytes (or roughly 77.5 MB) per image, assuming a high estimate for cancerous cases with 0.1% of pixels being non-zero.

CODE IS LINKED [HERE](#):

**Q2**.
To simulate the images captured by the microscope and the dye sensor, I'll write Python code that reflects the characteristics and data structures chosen previously. The microscope images will be simulated using a run-length encoding (RLE) approach, and the dye sensor images will be represented using a sparse matrix format.

**Microscope Image Simulation**
For the microscope images, we'll generate a binary image with large blobs representing the parasites. To keep it simple yet somewhat realistic, we'll simulate a blob as a large rectangle of black pixels on a white background. We'll then apply RLE to this simulated image.

**Dye Sensor Image Simulation**
For the dye sensor images, we'll randomly place lit pixels within the boundary of the simulated parasite (blob) to represent the luminescent dye. There will also be a small number of

randomly placed lit pixels outside of this boundary to simulate dye leakage. We'll store these points in a sparse matrix format.

Let's start with the code for generating these simulated images. Note that due to the complexity of generating high-resolution images with Python in this environment, we'll work with a much smaller scale for demonstration purposes (e.g., 1000x1000 pixels instead of 100,000x100,000).

CODE IS LINKED [HERE](#):

**Q3**.
To determine whether a simulated parasite has cancer, we need to calculate the percentage of the area occupied by the dye within the boundary of the parasite. Recall that a parasite is deemed to have cancer if the total amount of dye detected in its body exceeds 10% of the area occupied by the parasite in the image.

To achieve this, we'll follow these steps:
1. Decode the RLE-encoded microscope image to determine the total area occupied by the parasite. This step is necessary because our determination of whether the parasite has cancer depends on comparing the lit dye area to the total parasite area.
2. Count the number of lit pixels within the boundary of the parasite using the list of lit pixels generated for the dye sensor image. We already have this list from the simulation, but we need to ensure these pixels are within the parasite's area.
3. Compute the percentage of the area lit by dye and compare it to the 10% threshold to determine if the parasite has cancer.

CODE IS LINKED [HERE](#):

**Q4**.
To optimize the function for determining if a parasite has cancer, especially for high-resolution images (100,000 x 100,000 pixels), we need to focus on minimizing computational overhead and leveraging efficient data processing libraries like NumPy. Given the simplicity of the simulation, the primary bottleneck in a real-world scenario would likely come from the decoding process of the RLE data and the computation over large arrays. Let's address these concerns with a more optimized approach:

1. Avoid Decoding the Entire Image: Instead of decoding the whole RLE to check the dye percentage, we can directly work with the RLE data where possible, or

optimize the sparse representation of the dye sensor image to quickly intersect with the parasite's area.

2. Use Efficient Data Structures: For the sparse representation of the dye sensor data, ensure that we're using a structure that allows for rapid querying and intersection with the decoded blob area.

CODE IS LINKED [HERE](#):

## Q5.

Dealing with images of size 100,000x100,000 pixels presents significant challenges in terms of both storage and processing. The traditional compression techniques and runtime optimizations discussed earlier might not be sufficient or feasible at this scale. Below are some additional compression techniques that could be considered, along with their potential impacts on runtime and storage:

**For Microscope Images:**

**Quadtree Compression:** This technique is particularly effective for images with large areas of uniform color, as it recursively divides the image into quadrants until each quadrant can be represented as a single color. This could significantly reduce the storage for the binary microscope images.

- Runtime Impact: Compression might be slower due to the recursive division, but decompression can be very fast, especially for accessing specific image regions.

**Bitmap Indexing with Run-Length Encoding (RLE)**: This approach involves using a bitmap to represent the presence (or absence) of the parasite in each pixel, combined with RLE for further compression. It's a slight variation on the initial proposal with potentially better compression for specific patterns.

- Runtime Impact: Similar to RLE, but with potentially better compression and thus slightly increased computation for encoding/decoding.

**For Dye Sensor Images:**

**Sparse Matrix Formats (Compressed Sparse Row - CSR or Compressed Sparse Column - CSC)**: Given the expectation of sparse lit areas, CSR or CSC formats could efficiently store non-zero elements with less overhead than COO (Coordinate List), especially for operations like matrix multiplication or dot products, which might not be directly applicable here but are useful for image processing.

- Runtime Impact: These formats can offer faster access and operations on the sparse matrices compared to COO, especially if the matrix operations are optimized for these formats.

**Point Cloud Compression**: Considering the lit areas can be sparse and scattered, treating them as a point cloud and applying point cloud compression algorithms might be beneficial. This approach is more common in 3D data but could be adapted for 2D.

- Runtime Impact: Compression and decompression times can vary significantly based on the algorithm's complexity and the data's sparsity.

**Measuring Runtime and Storage Costs**

Unfortunately, performing actual runtime and storage cost measurements on my setup for 100,000x100,000 pixel images is not feasible due to the limitations of the environment in which this model operates. Such experiments would require significant computational resources (e.g., high RAM for in-memory processing and substantial disk space for storage) and specialized software or libraries capable of efficiently handling data of this magnitude. For accurate measurements, we would typically use profiling tools available in our programming environment (such as Python's cProfile for runtime profiling and custom logging for memory usage) and conduct the experiments on a machine with sufficient resources. The approach would involve:

- Implementing the compression technique for both types of images.
- Creating or obtaining sample images of the full 100,000x100,000 pixel resolution that are representative of the expected real data.
- Running the compression and decompression processes, measuring the time taken for each operation, and the storage size before and after compression.

To effectively handle and analyze the substantial amount of data we're working with, particularly our 100,000x100,000 pixel images, we'll need to leverage our local machines or cloud computing resources that are equipped for large-scale data processing. Here's our game plan for measuring the runtime of our compression operations on these large images:

1. **Image Preparation**

First up, we're going to generate binary images that simulate the output from our microscope and dye sensor. Given the immense size of these images, it's crucial we use a library like NumPy, which is designed to efficiently manage large arrays. This step will set us up with a solid foundation for our subsequent compression tasks.

2. **Compression Implementation**

Next, we'll dive into coding our compression algorithms. The key here is to focus on both memory efficiency and processing speed. Since we're dealing with images of a massive scale, even minor optimizations can lead to significant improvements in performance. Let's ensure our code is as streamlined and efficient as possible.

3. **Runtime Measurement**

Finally, we'll measure the execution time of our compression operations. There are a couple of ways we can approach this:

4. **Using the time module:** This method involves simply wrapping our compression operation with start and end time captures. It's straightforward and gives us a quick snapshot of the operation's duration.

5. **Employing cProfile for detailed analysis**: For those of us looking for a deeper dive into the execution time of various parts of our code, cProfile is the way to go. It can provide us with detailed insights into where we might further optimize our processes.

Here's a simplified example using NumPy and time to measure how long it takes to create the image and perform a basic operation, assuming we're working with a smaller, more manageable image size for demonstration: CODE IS LINKED [HERE](#):

**Q6**.
To tackle this challenge, I primarily leaned on Python, especially its NumPy library, for managing and manipulating the large image arrays. NumPy is a cornerstone in Python data science and offers extensive support for large, multi-dimensional arrays and matrices, which was crucial for simulating and processing our 100,000x100,000 pixel images efficiently. This choice was driven by the need to handle complex operations and data structures with both speed and efficiency, ensuring that our simulations and analyses remained practical despite the vast scale of the data involved. In addition to Python and NumPy, I also integrated insights and methodologies derived from ChatGPT, particularly for conceptual understanding and optimizing algorithms. ChatGPT served as a virtual collaborator, providing on-the-fly suggestions, code snippets, and theoretical insights that were instrumental in refining our approach to data compression and analysis. The interaction with ChatGPT allowed me to quickly explore various strategies and troubleshoot complex problems, significantly accelerating the development process.