The following scalability considerations and techniques are suggested to guarantee that the OCR system can effectively handle large-scale cadastral map datasets:

1. **Flow Design:**
   a. The existing pipeline processes each image sequentially:
      i. Read image
      ii. Pre-process image
      iii. Apply OCR
      iv. Post-process text
      v. Write numbers and decimals into CSV
   b. Processing images sequentially increases latency as the dataset size increases.
   c. To ease the burden on resources, images can be processed in batches.
2. **Parallel Processing:**
   a. Image processing can be parallelised by using Python's multiprocessing or "concurrent.futures" module.
   b. Serverless functions can be used to install the pipeline on cloud platforms like Amazon Web Services (AWS).
3. **Efficient Data Handling Techniques:**
   a. Photographs can be uploaded and processed in lightweight data forms, such as compressed images.
   b. For quick access, images can be kept in a distributed file system like HDFS or AWS S3.
   c. OCR models and other frequently used data can be routinely cached.