



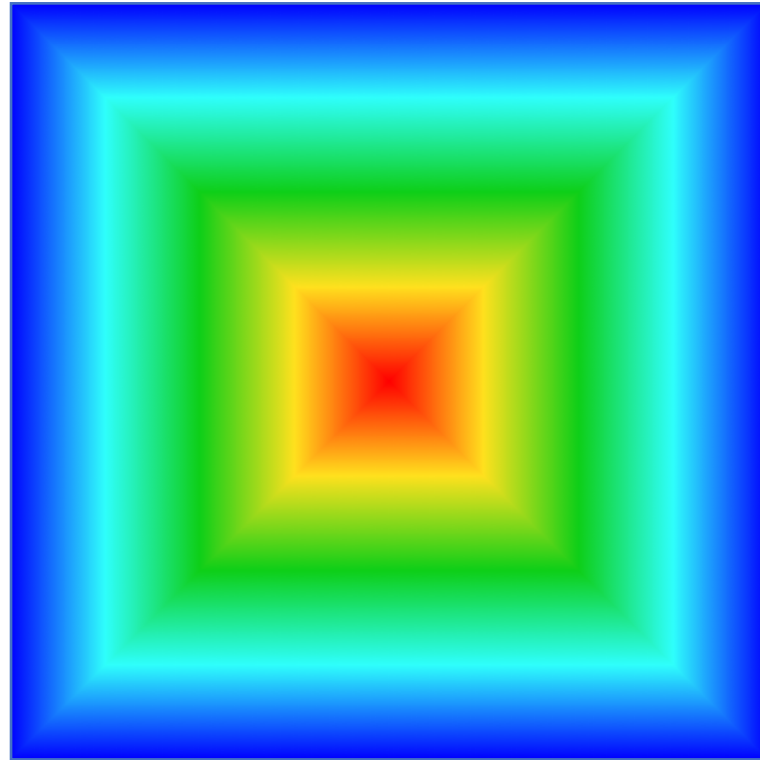
The Abdus Salam  
International Centre  
for Theoretical Physics

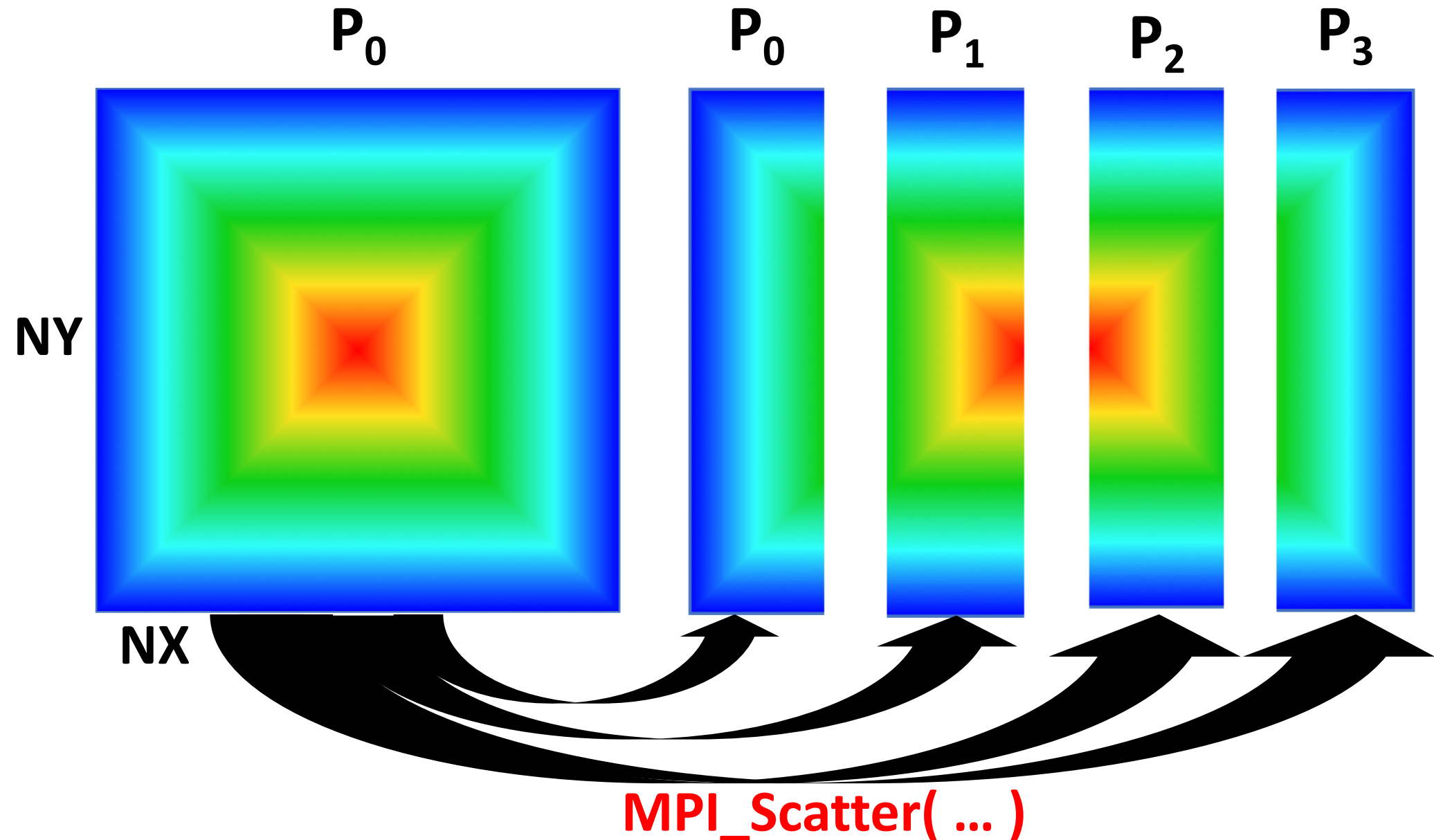


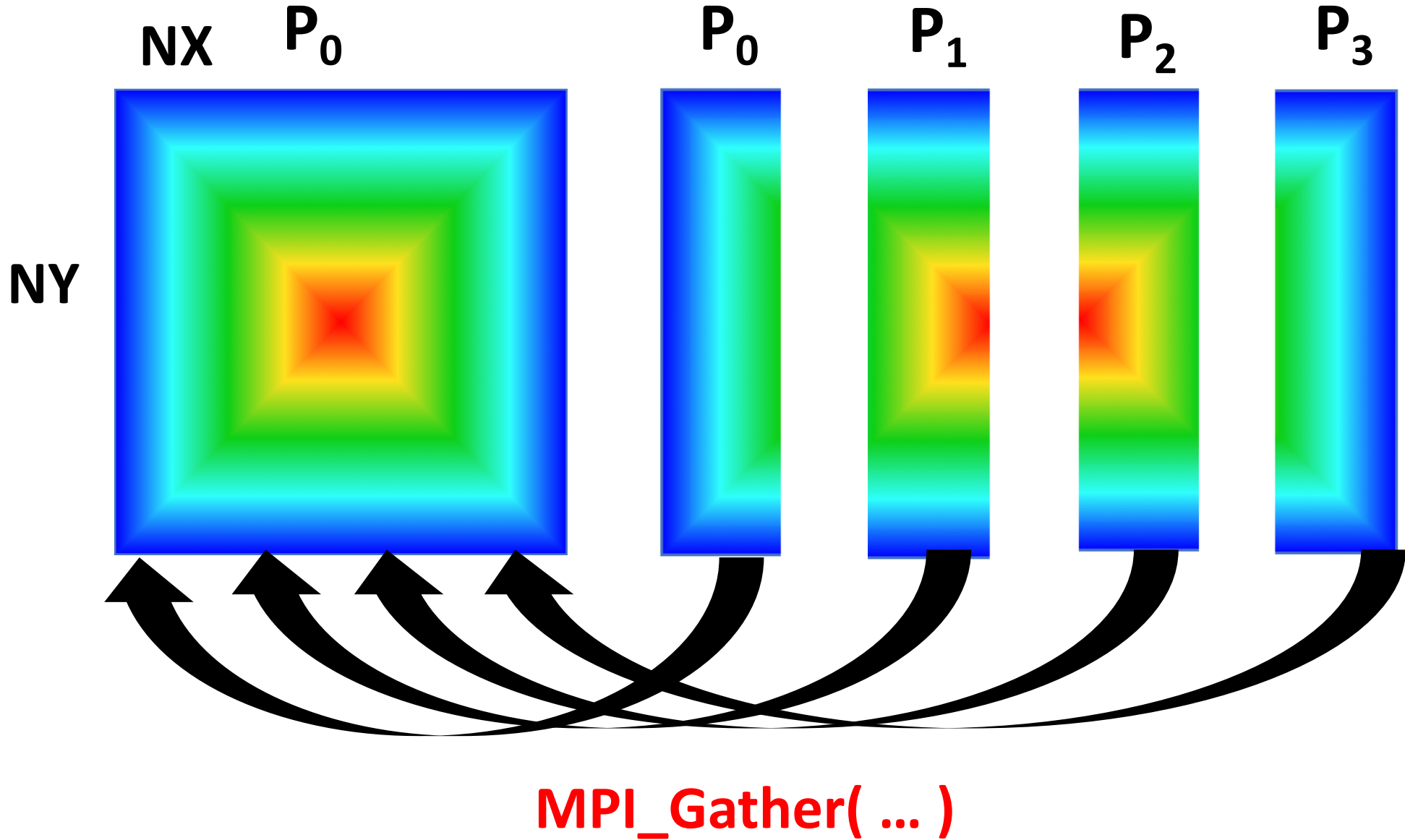
# Foundation of Parallel I/O

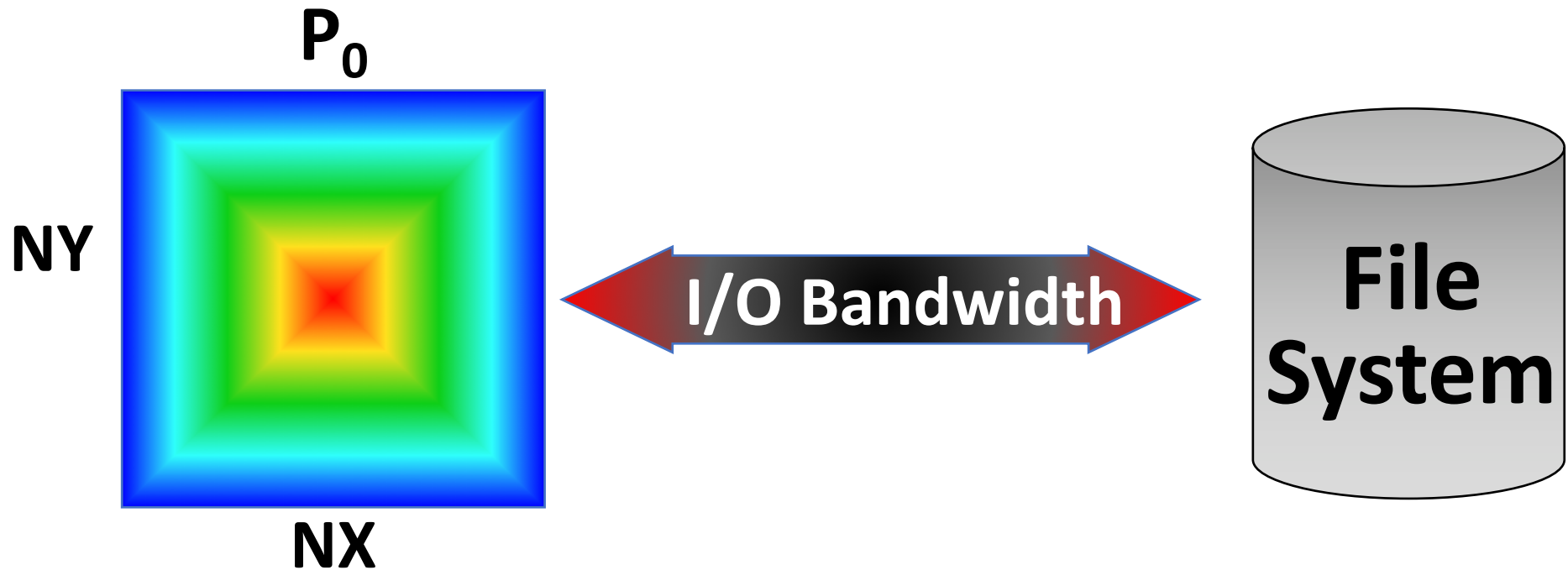
**Ivan Girotto – [igirotto@ictp.it](mailto:igirotto@ictp.it)**

International Centre for Theoretical Physics (ICTP)

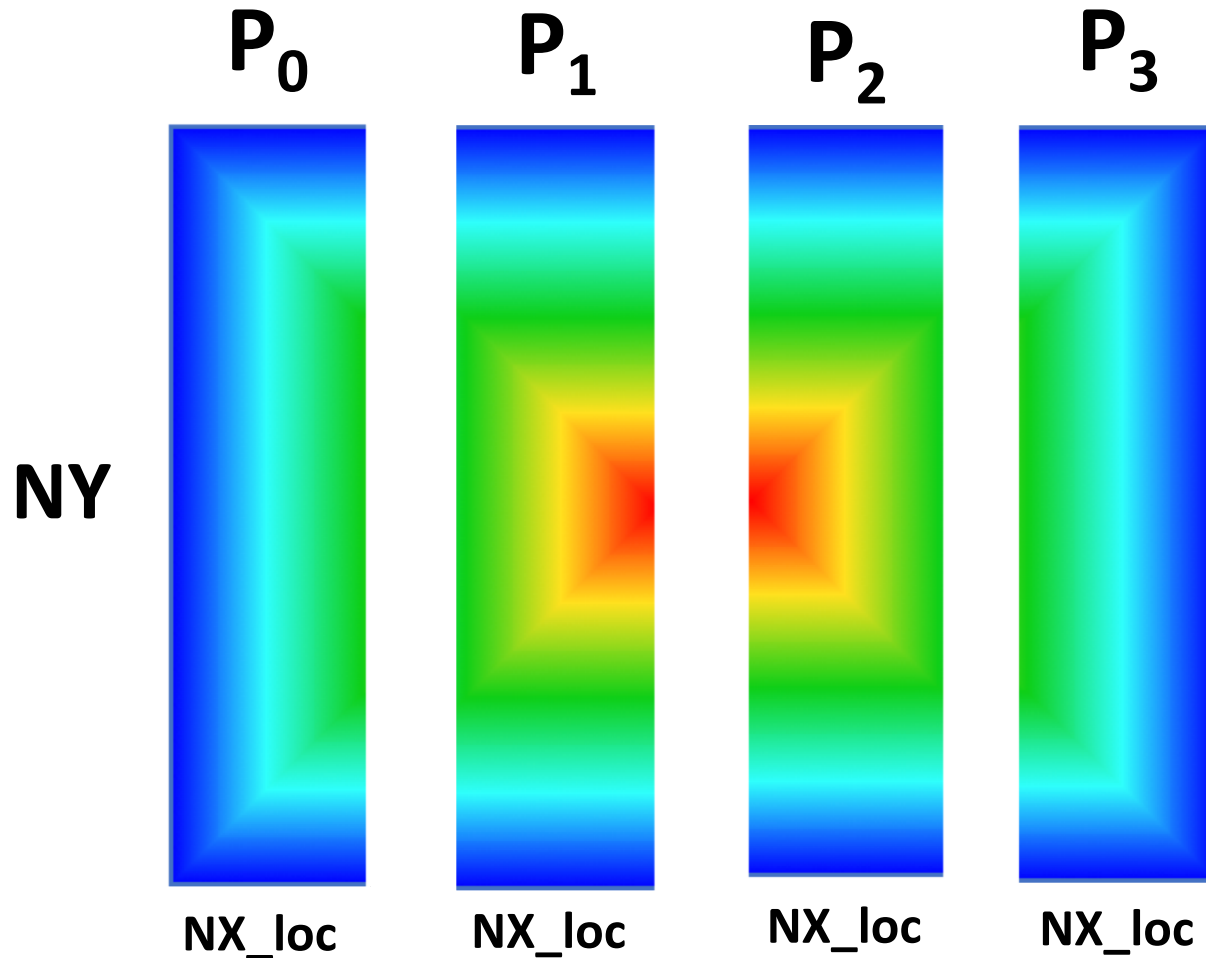






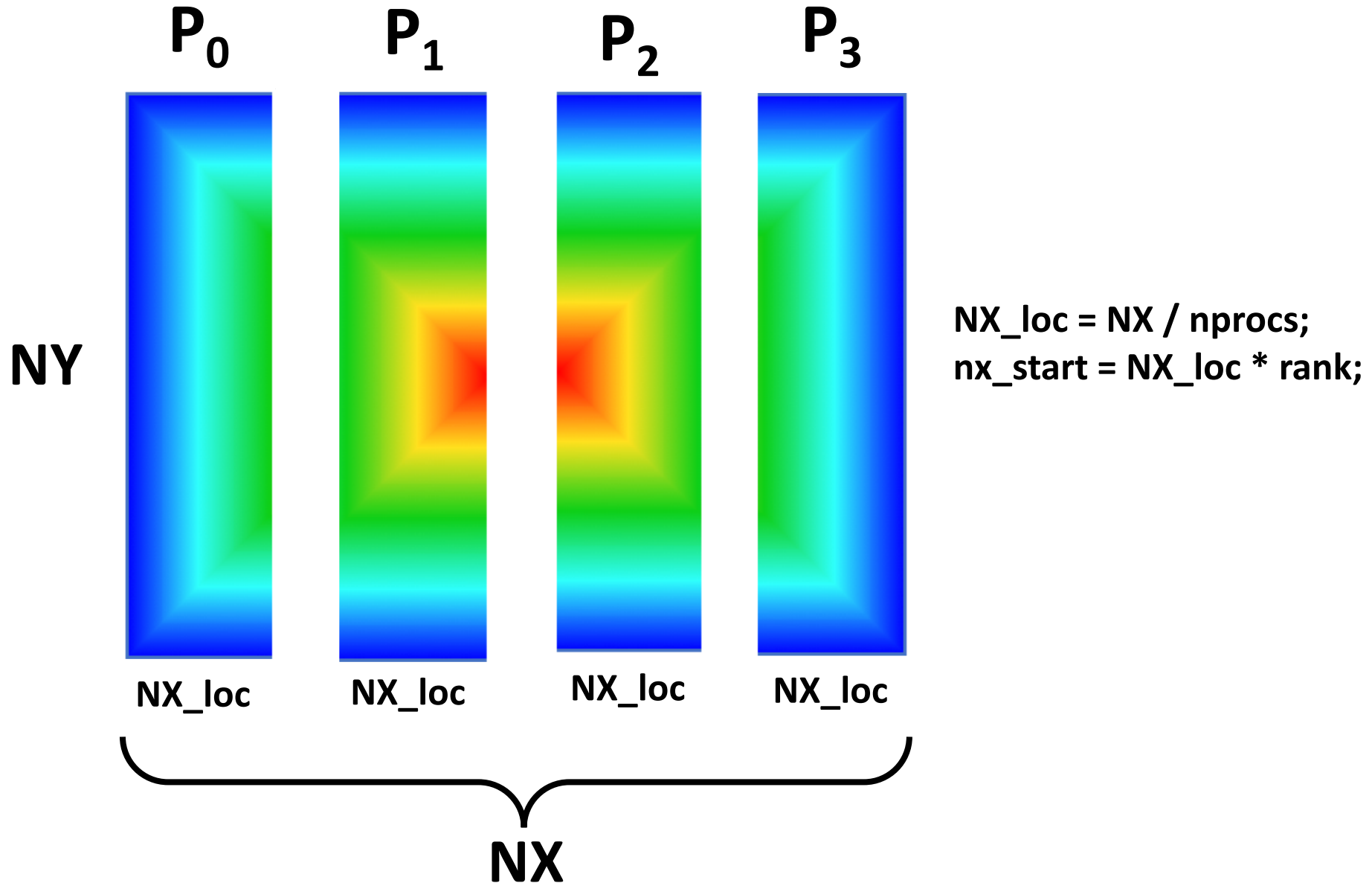


Data are finally saved in the disk like on the model.cpp code by using HDF5. By collecting all data on a single process we are implementing a strong data size scaling limitation!

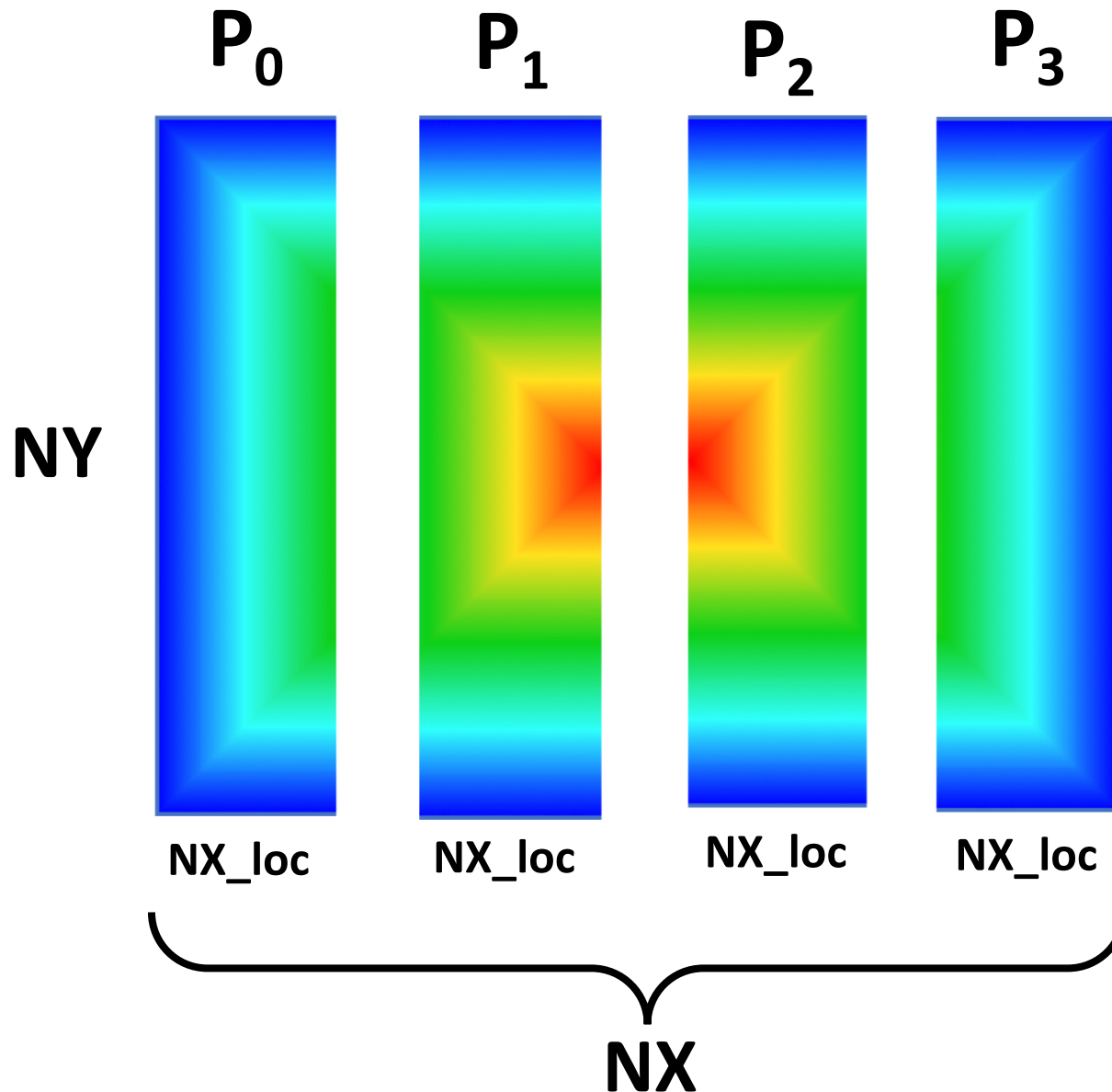


We now want to keep data distributed into memory with never collecting data on a single processor to achieve massively data size scaling!!

The lucky case:  $NX \% nprocs == 0$  !!!



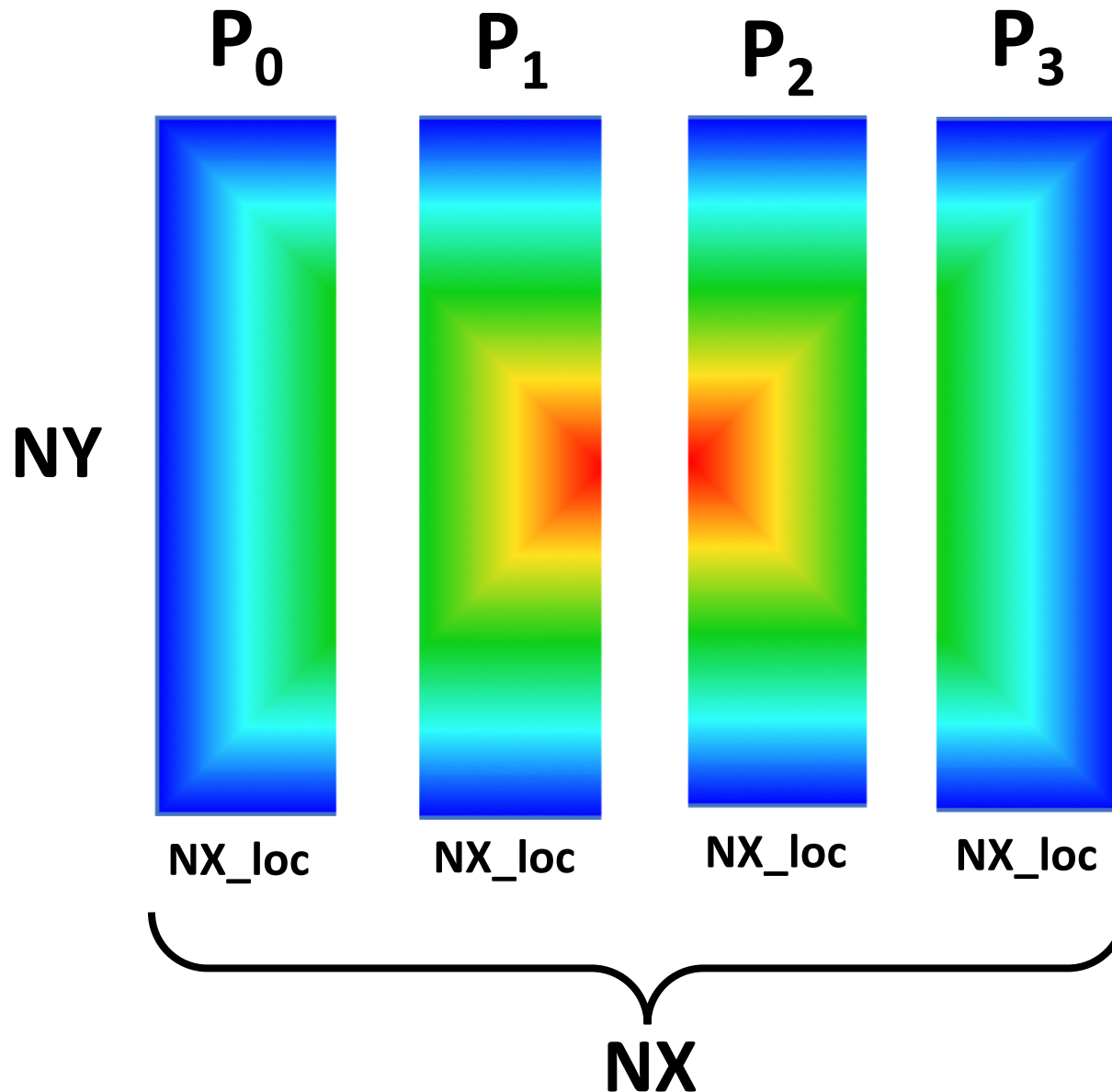
The unlucky case:  $NX \% nprocs \neq 0$  !!!



```
NX_loc = NX / nprocs;  
rest = NX % nprocs;  
nx_start = NX_loc * rank;  
If ( me == nprocs - 1 ) NX_loc += rest;
```



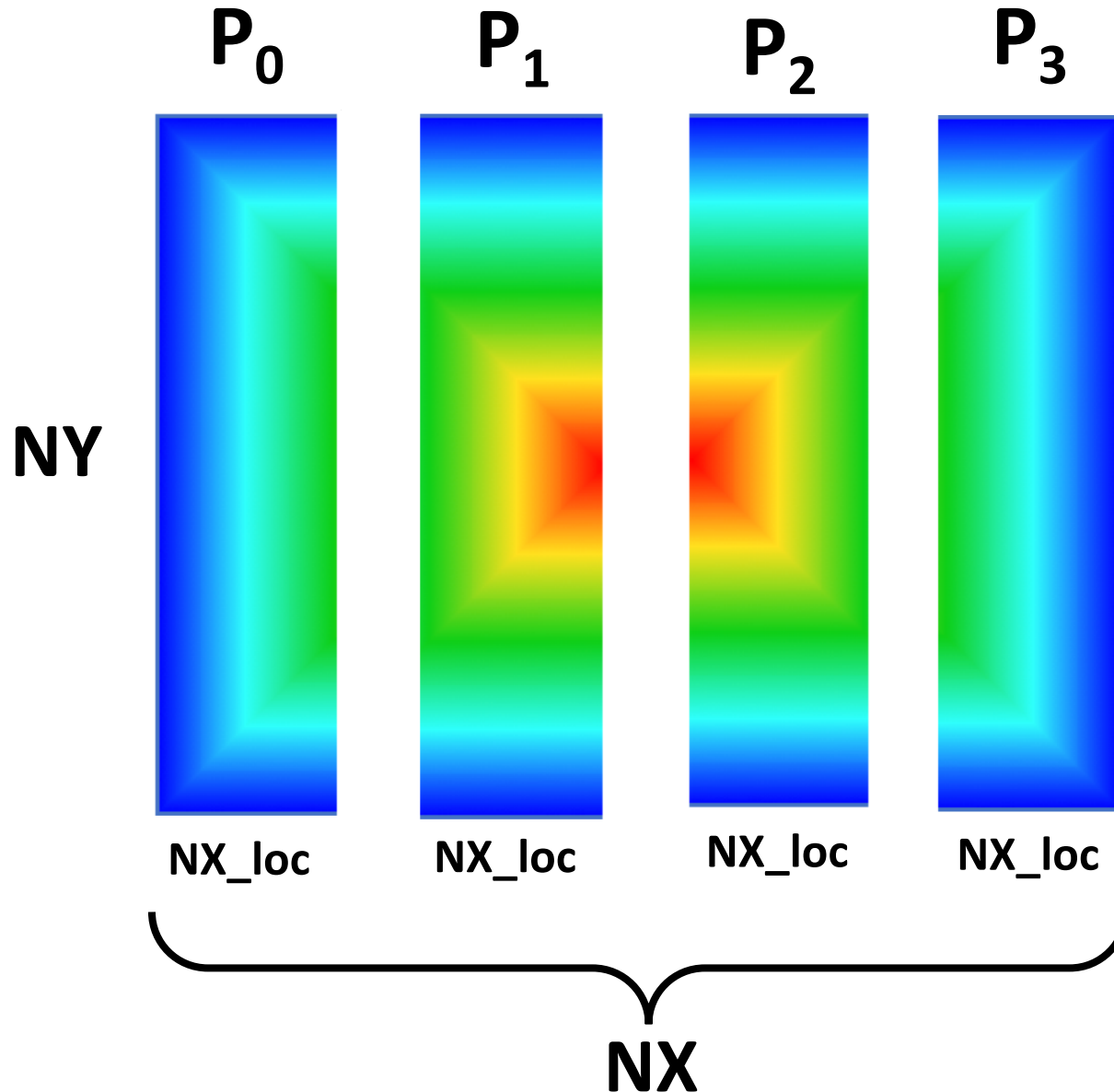
The unlucky case:  $NX \% nprocs \neq 0$  !!!



```
NX_loc = NX / nprocs;  
rest = NX % nprocs;  
nx_start = NX_loc * rank;  
If ( me == nprocs - 1 ) NX_loc += rest;
```

Example:  
 $NX = 1400000$ ;  
 $nprocs = 480$ ;

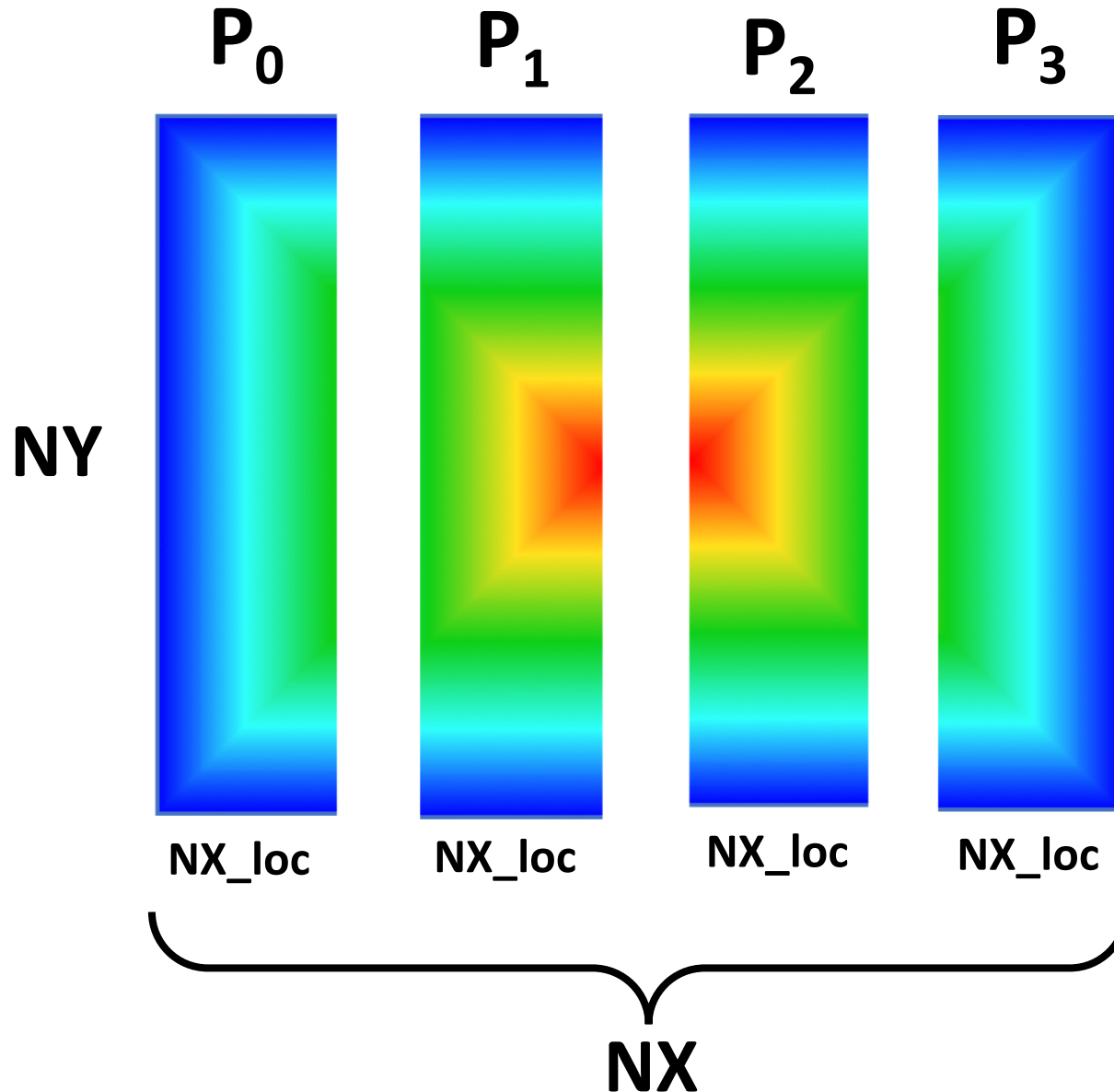
The unlucky case:  $NX \% nprocs \neq 0$  !!!



```
NX_loc = NX / nprocs;  
rest = NX % nprocs;  
nx_start = NX_loc * rank;  
If ( me == nprocs - 1 ) NX_loc += rest;
```

Example:  
 $NX = 1400000$ ;  
 $nprocs = 480$ ;  
 $NX\_loc = 291$

The unlucky case:  $NX \% nprocs \neq 0$  !!!



```
NX_loc = NX / nprocs;  
rest = NX % nprocs;  
nx_start = NX_loc * rank;  
If ( me == nprocs - 1 ) NX_loc += rest;
```

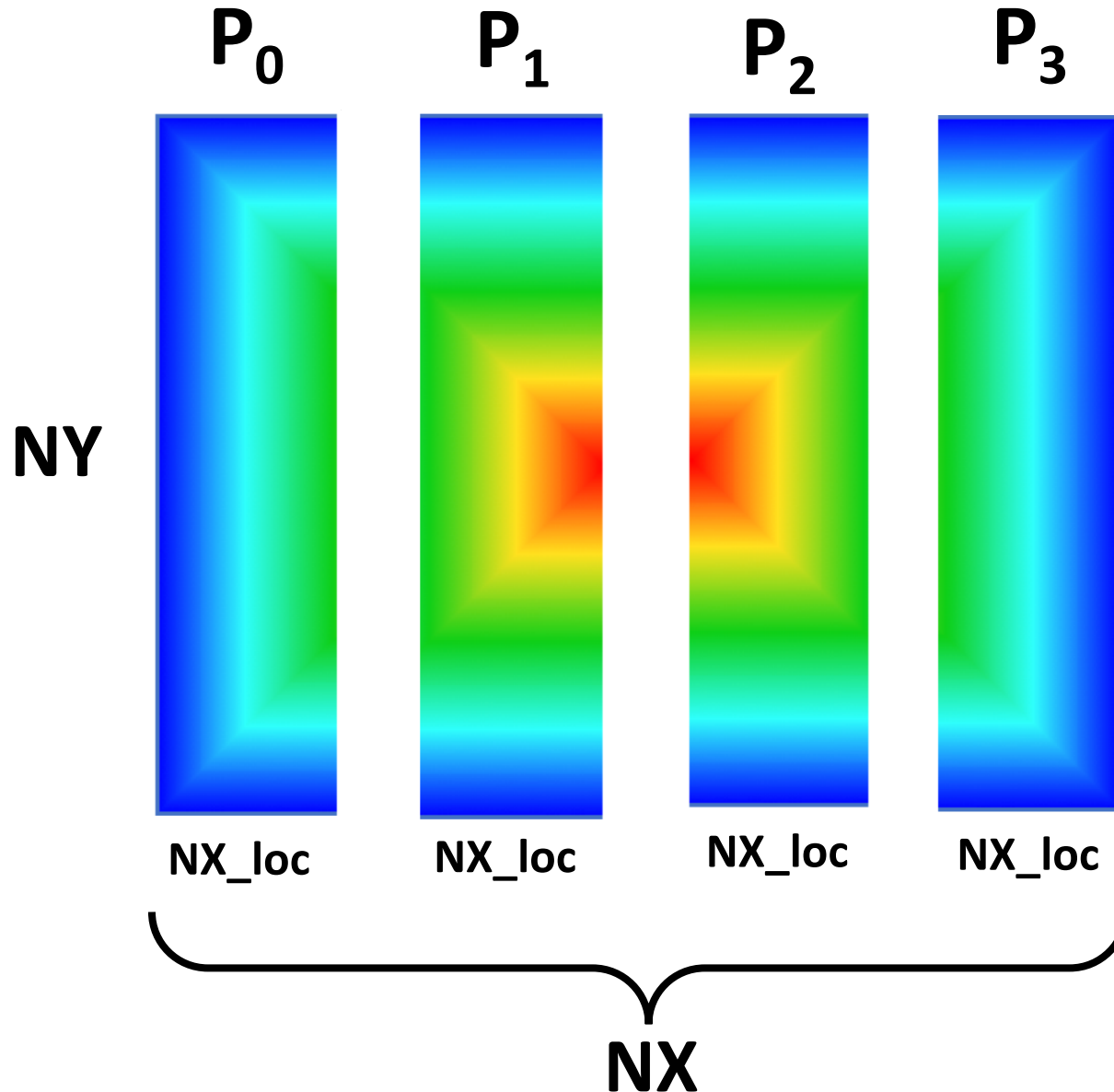
Example:

```
NX = 1400000;  
nprocs = 480;  
NX_loc = 291;  
rest = 320;
```

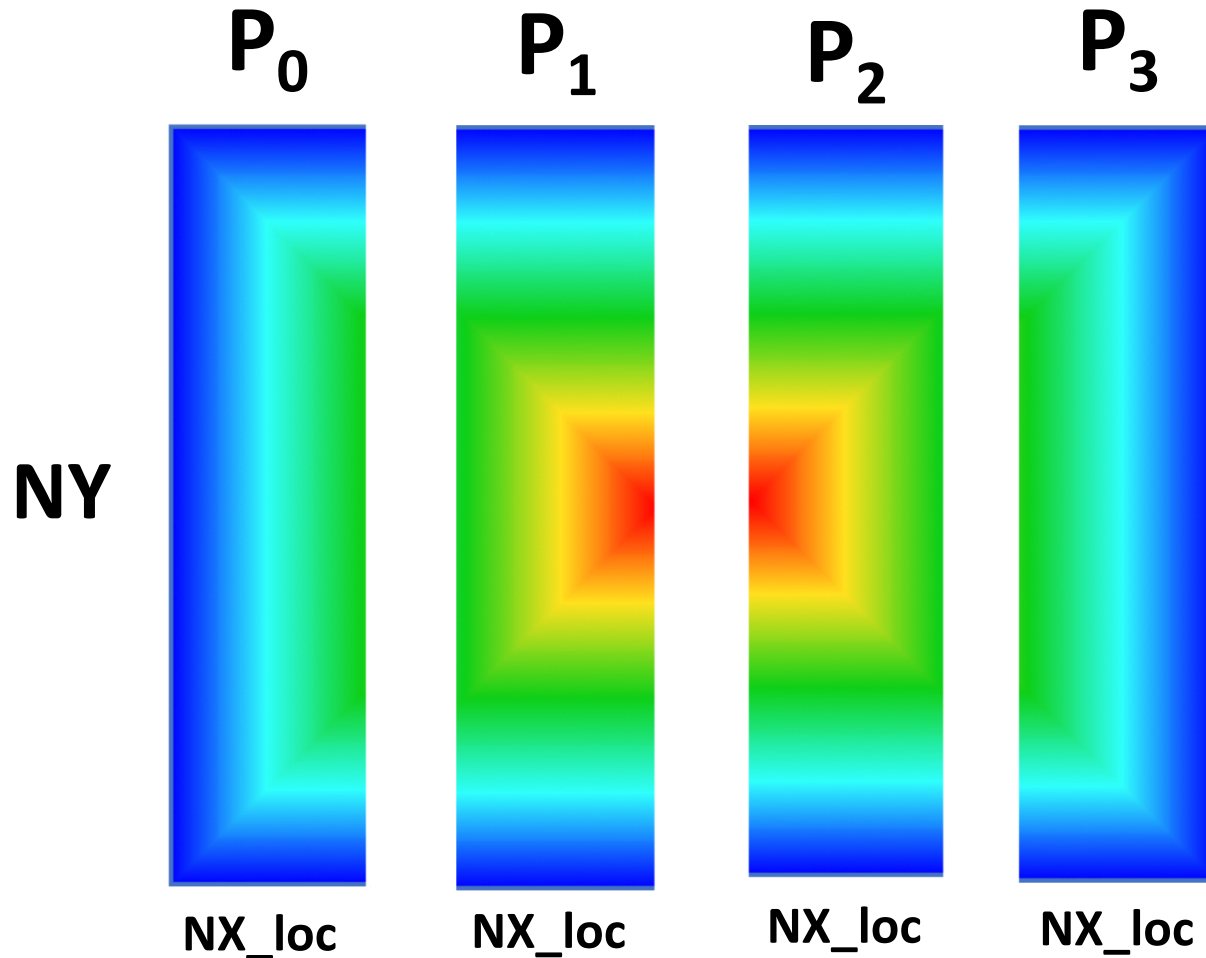
$NX\_loc (nprocs - 1) = 611!!$

**NO WAY!!!!**

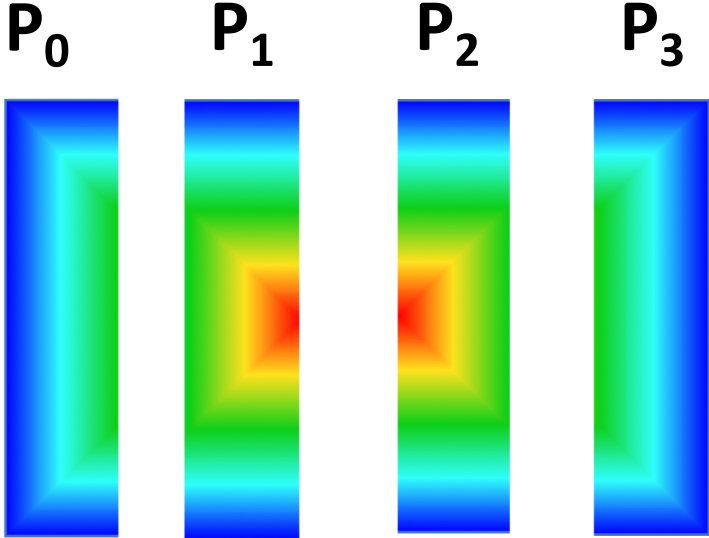
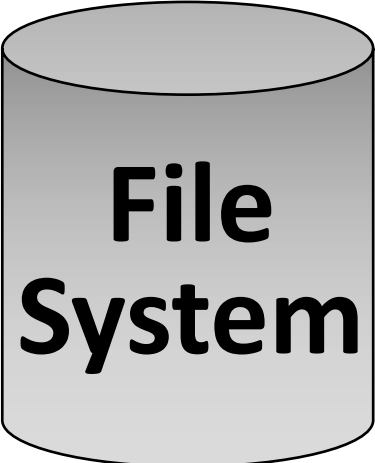
The unlucky case:  $NX \% nprocs \neq 0$  !!!



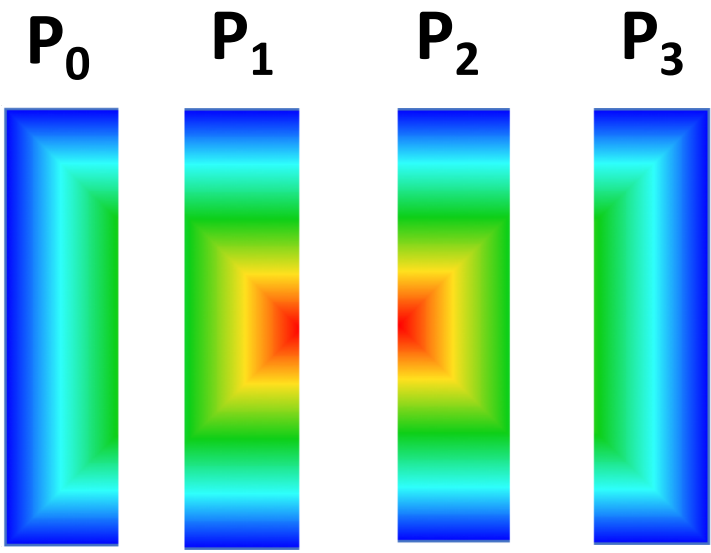
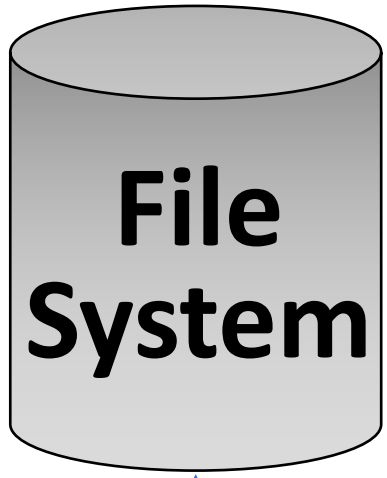
```
NX_loc = NX / nprocs;  
rest = NX % nprocs;  
If ( me < rest ) NX_loc += 1;  
nx_start = NX_loc * rank;  
If ( me >= rest ) nx_start += rest;
```



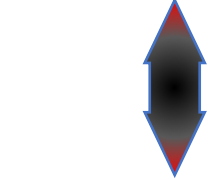
We now want to save data on the disk but we do not have a parallel file system!

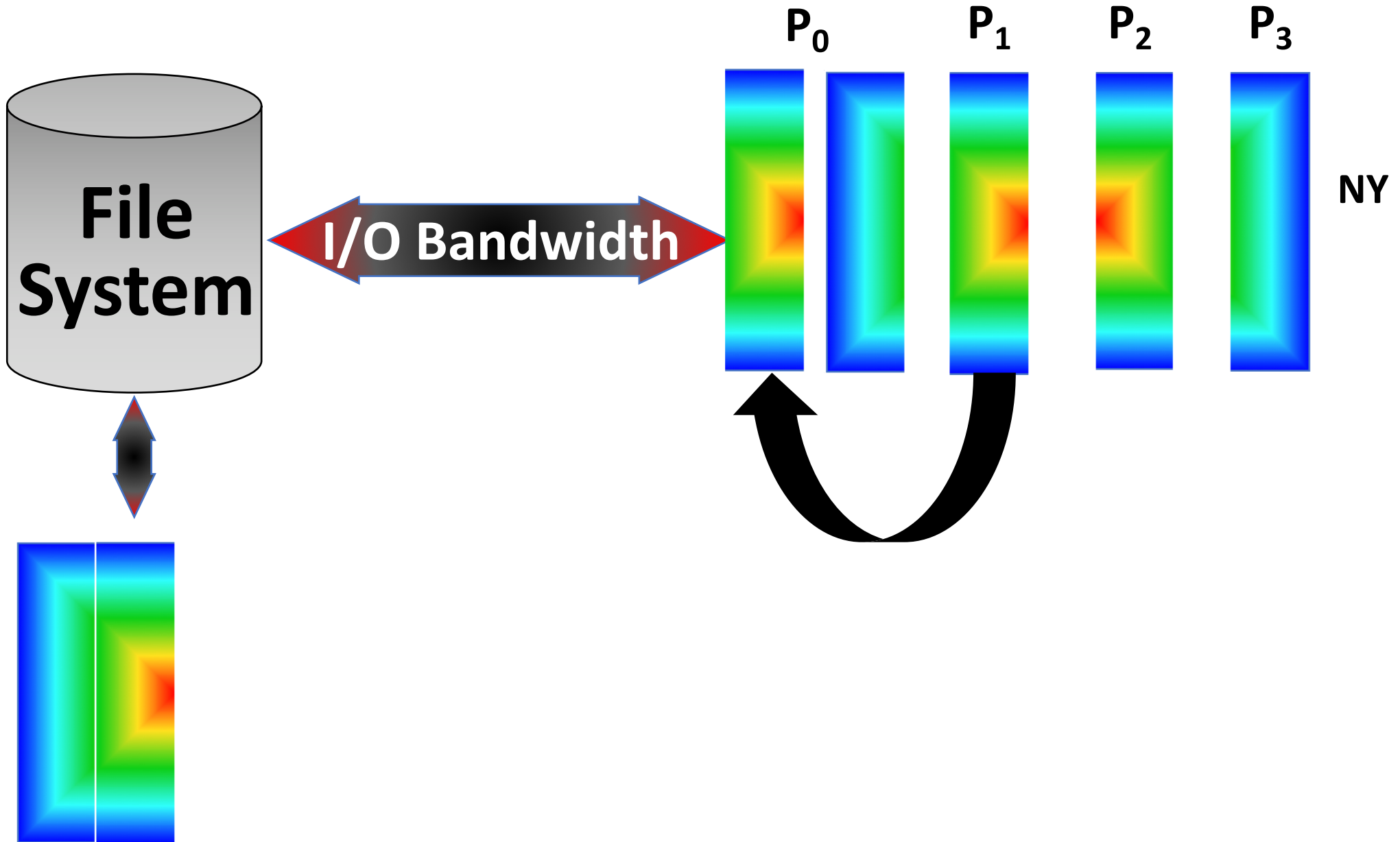


NY

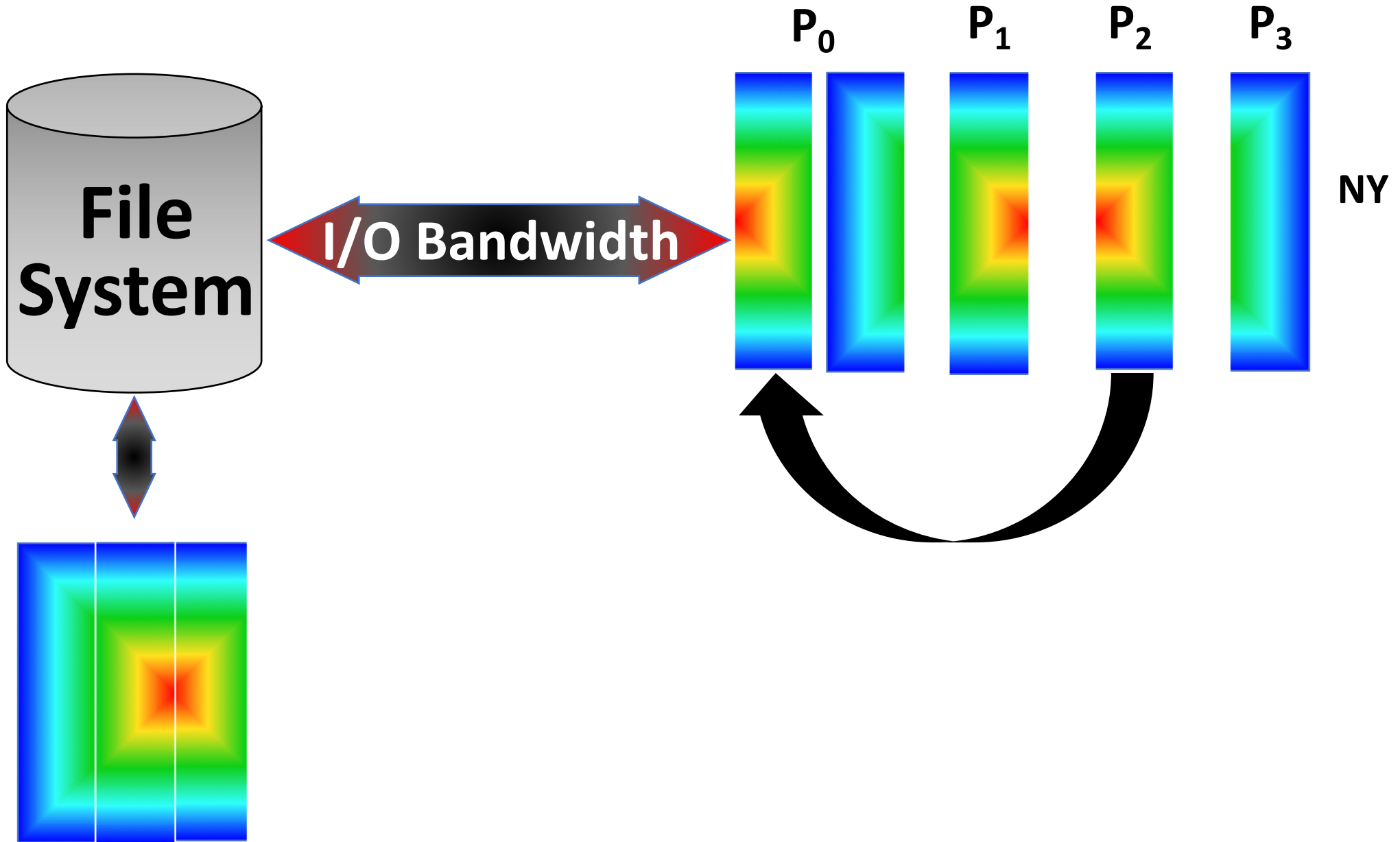


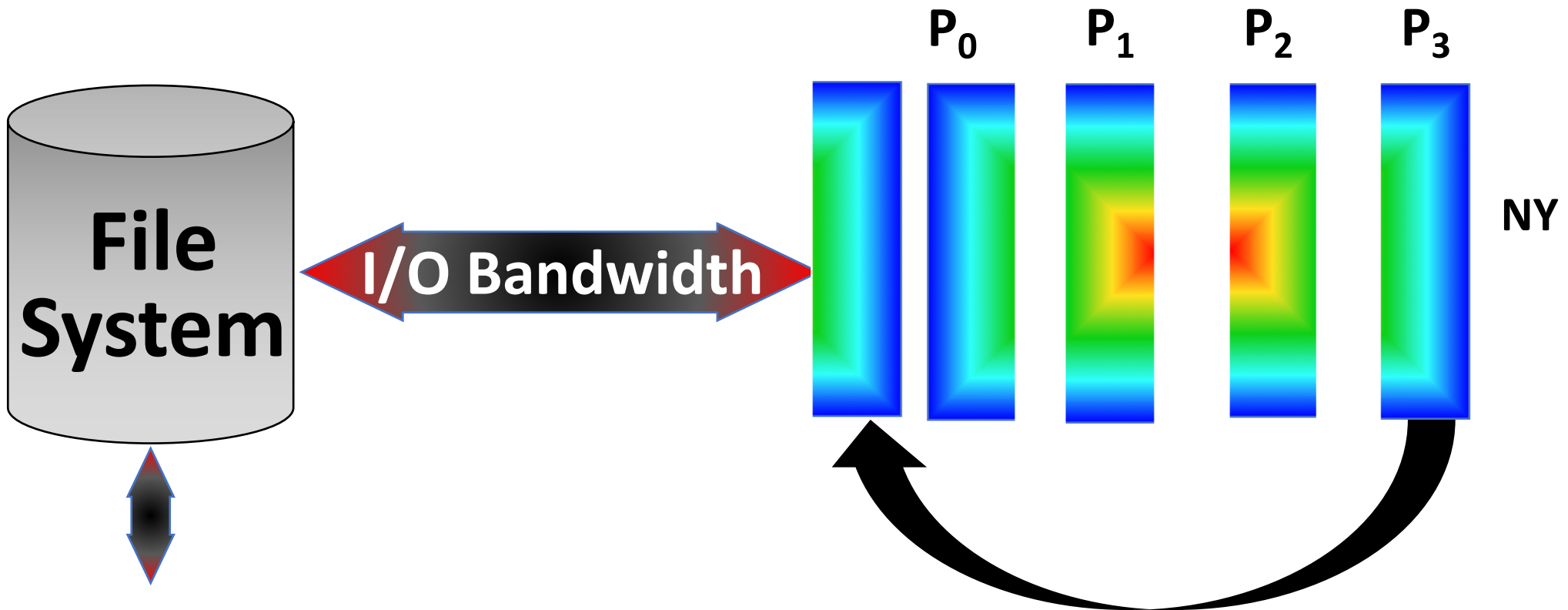
NY



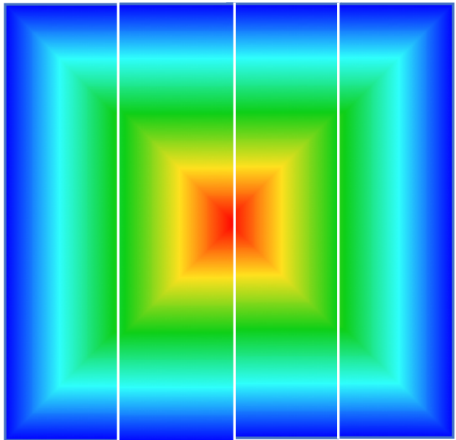


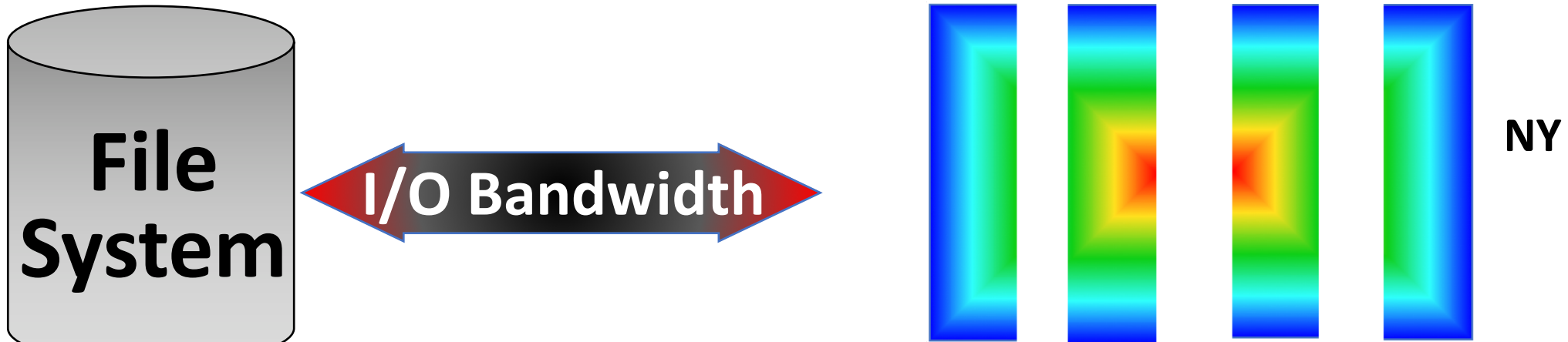




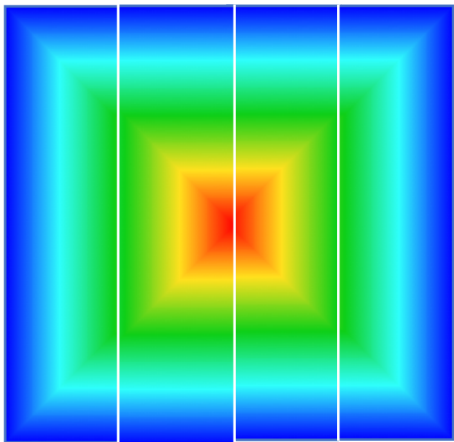


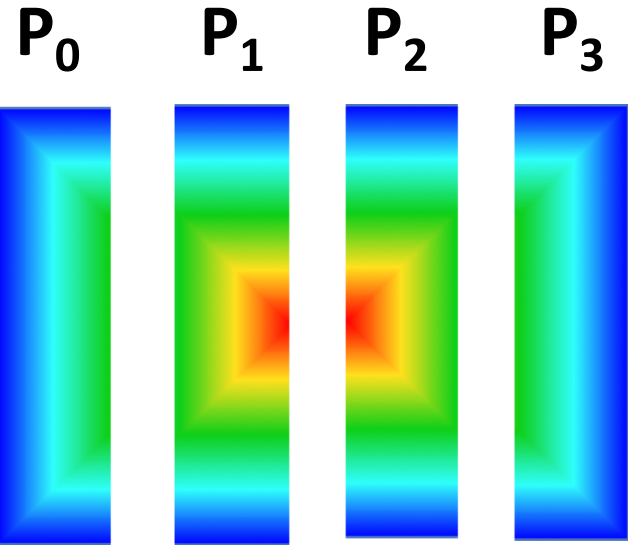
This strategy works in most cases but it might bring together a significant unbalance on the I/O process!!





- A strategy where only one process (node) is dedicated to perform operations of I/O
- On a MPI based parallelization it requires the creation of additional communicators
- The I/O process should be isolated on a single node if allocating much larger memory than others
- For higher scalability multiple I/O can be implemented
  - Non trivial MPI execution model (but possible)



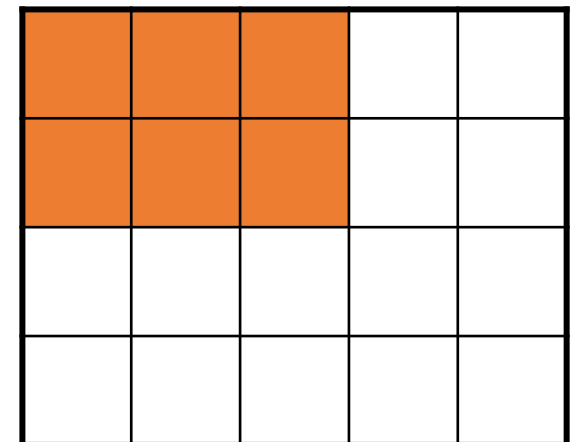


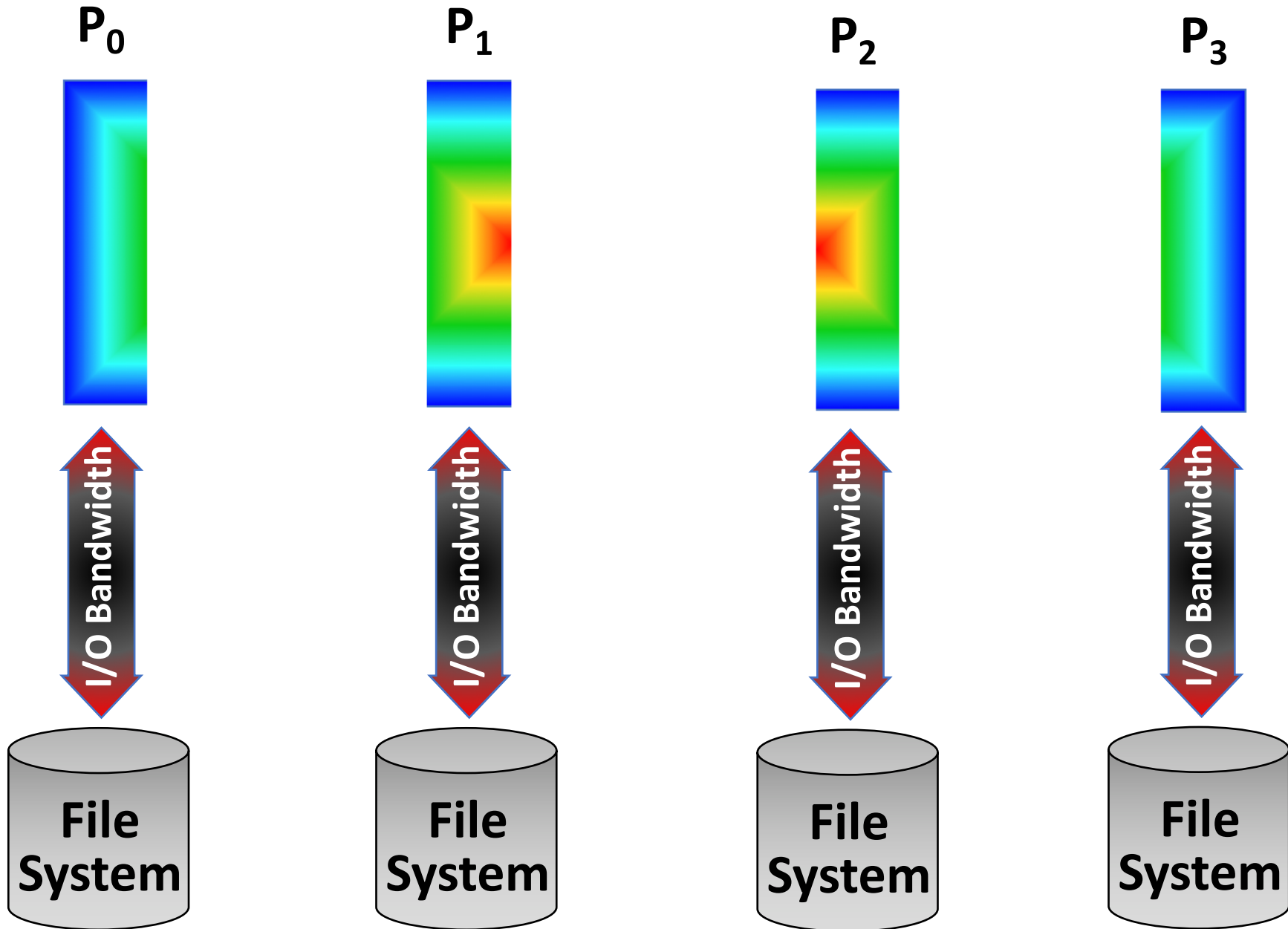
- Parallel Reading is always possible
- No concurrent access on data
- All read all => memory scaling issue!!
- It is really efficient when all read a chunk of data creating a distributed data schema

- Parallel Reading is always possible
- No concurrent access on data
- All read all => memory scaling issue!!
- It is really efficient when all read a chunk of data creating a distributed data schema

## MPI Data Types!!

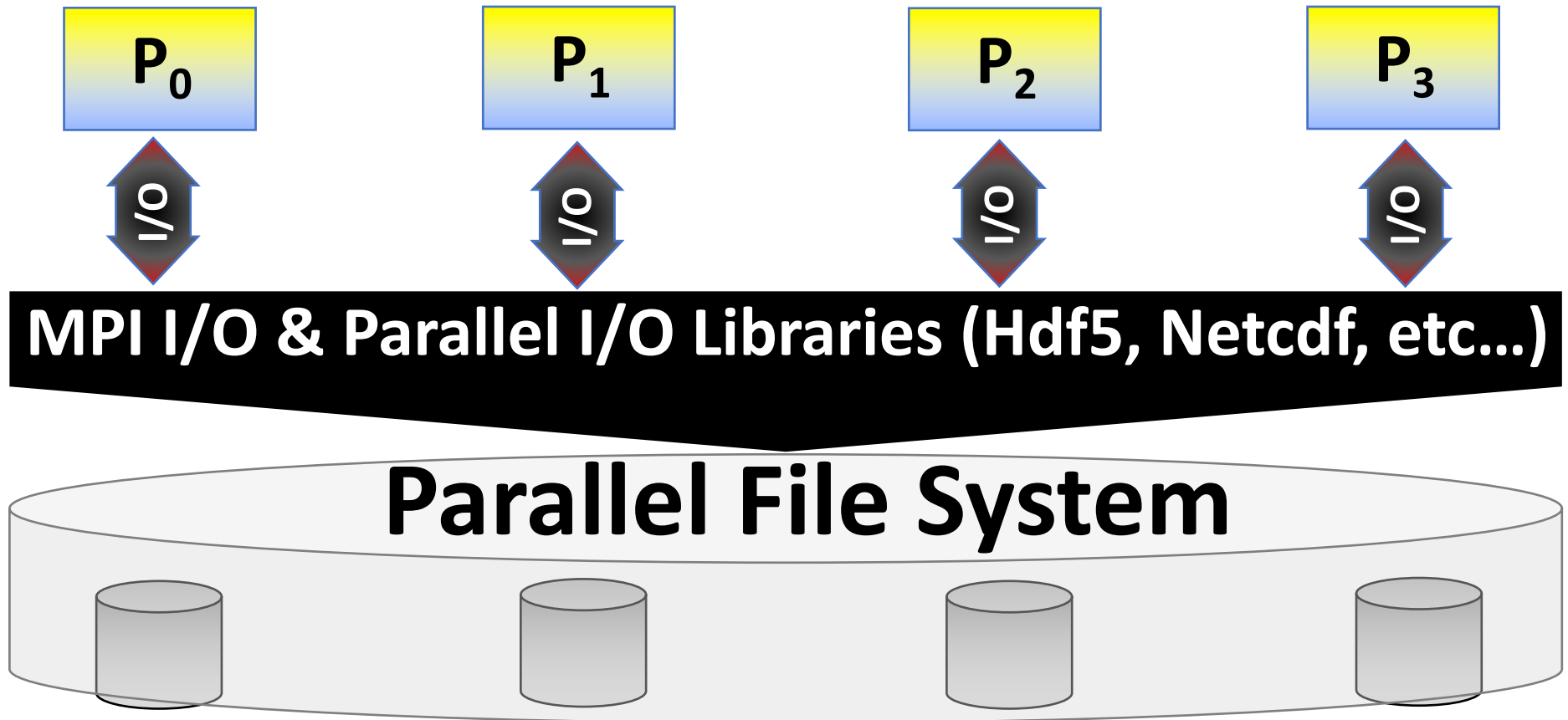
`MPI_Type_vector(2, 3, 5, MPI_INT, MPI_SUBMAT)`





- Every process writes its own file
- Easy to implement as long as every file is identified with a unique file name
- Makes really difficult to restart with a different number of processes
- Might create enormous problem at post processing
  - If we need to recreate the whole data-set we need to read the data again and do what we were supposed to do initially when data were first moving from memory to disk
  - Possible creation of an enormous number of files. You could create a problem not only to your self but to the whole infrastructure

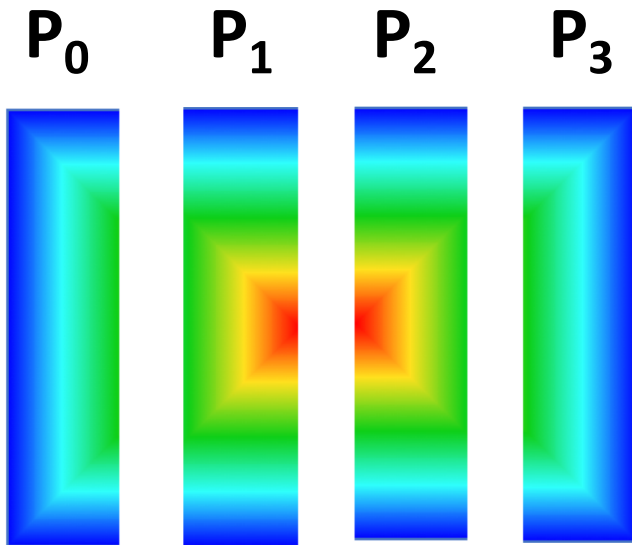




- The Parallel File system provide the utility to write a file from distributed processes on the storage system
- The intermediate level is essential to avoid learning how to handle concurrent access to shared file pointers :-o

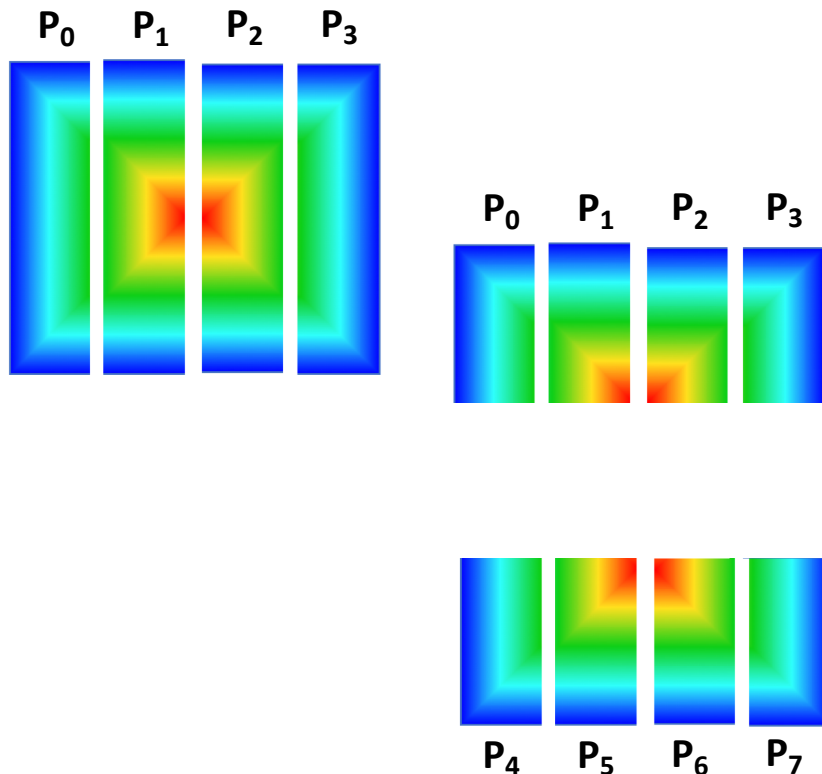
# Parallel I/O: MPI I/O

- Use of MPI routines for parallel I/O
- Interfaces for handling shared access on the same file
- Developers must implement a way to store data accordingly to a given data layout => MPI data types helps!

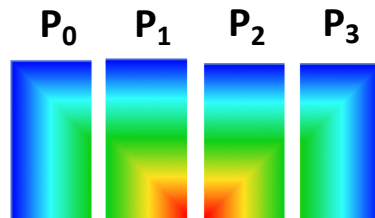
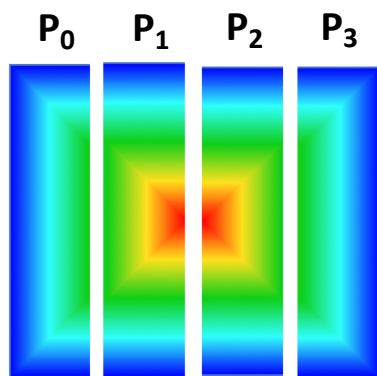


# Parallel I/O: MPI I/O

- Use of MPI routines for parallel I/O
- Interfaces for handling shared access on the same file
- Developers must implement a way to store data accordingly to a given data layout => MPI data types helps!



- Use of MPI routines for parallel I/O
- Interfaces for handling shared access on the same file
- Developers must implement a way to store data accordingly to a given data layout => MPI data types helps!

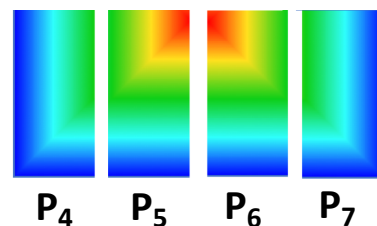


a <sub>11</sub>	a <sub>12</sub>	a <sub>13</sub>	a <sub>14</sub>	a <sub>15</sub>	a <sub>16</sub>	a <sub>17</sub>	a <sub>18</sub>	a <sub>19</sub>
a <sub>21</sub>	a <sub>22</sub>	a <sub>23</sub>	a <sub>24</sub>	a <sub>25</sub>	a <sub>26</sub>	a <sub>27</sub>	a <sub>28</sub>	a <sub>29</sub>
a <sub>31</sub>	a <sub>32</sub>	a <sub>33</sub>	a <sub>34</sub>	a <sub>35</sub>	a <sub>36</sub>	a <sub>37</sub>	a <sub>38</sub>	a <sub>39</sub>
a <sub>41</sub>	a <sub>42</sub>	a <sub>43</sub>	a <sub>44</sub>	a <sub>45</sub>	a <sub>46</sub>	a <sub>47</sub>	a <sub>48</sub>	a <sub>49</sub>
a <sub>51</sub>	a <sub>52</sub>	a <sub>53</sub>	a <sub>54</sub>	a <sub>55</sub>	a <sub>56</sub>	a <sub>57</sub>	a <sub>58</sub>	a <sub>59</sub>
a <sub>61</sub>	a <sub>62</sub>	a <sub>63</sub>	a <sub>64</sub>	a <sub>65</sub>	a <sub>66</sub>	a <sub>67</sub>	a <sub>68</sub>	a <sub>69</sub>
a <sub>71</sub>	a <sub>72</sub>	a <sub>73</sub>	a <sub>74</sub>	a <sub>75</sub>	a <sub>76</sub>	a <sub>77</sub>	a <sub>78</sub>	a <sub>79</sub>
a <sub>81</sub>	a <sub>82</sub>	a <sub>83</sub>	a <sub>84</sub>	a <sub>85</sub>	a <sub>86</sub>	a <sub>87</sub>	a <sub>88</sub>	a <sub>89</sub>
a <sub>91</sub>	a <sub>92</sub>	a <sub>93</sub>	a <sub>94</sub>	a <sub>95</sub>	a <sub>96</sub>	a <sub>97</sub>	a <sub>98</sub>	a <sub>99</sub>

Logical View (Matrix)

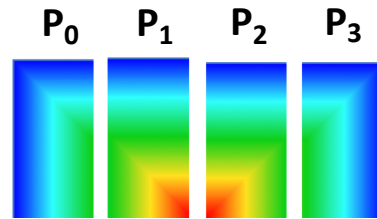
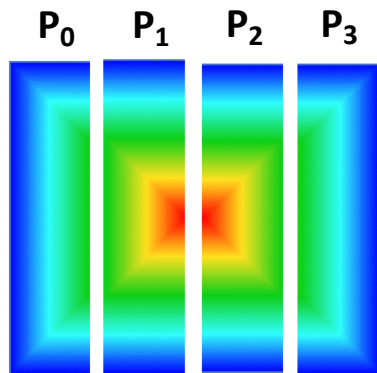
a <sub>11</sub>	a <sub>12</sub>	a <sub>17</sub>	a <sub>18</sub>	a <sub>13</sub>	a <sub>14</sub>	a <sub>19</sub>	a <sub>15</sub>	a <sub>16</sub>
a <sub>21</sub>	a <sub>22</sub>	a <sub>27</sub>	a <sub>28</sub>	a <sub>23</sub>	a <sub>24</sub>	a <sub>29</sub>	a <sub>25</sub>	a <sub>26</sub>
a <sub>51</sub>	a <sub>52</sub>	a <sub>57</sub>	a <sub>58</sub>	a <sub>53</sub>	a <sub>54</sub>	a <sub>59</sub>	a <sub>55</sub>	a <sub>56</sub>
a <sub>61</sub>	a <sub>62</sub>	a <sub>67</sub>	a <sub>68</sub>	a <sub>63</sub>	a <sub>64</sub>	a <sub>69</sub>	a <sub>65</sub>	a <sub>66</sub>
a <sub>91</sub>	a <sub>92</sub>	a <sub>97</sub>	a <sub>98</sub>	a <sub>93</sub>	a <sub>94</sub>	a <sub>99</sub>	a <sub>95</sub>	a <sub>96</sub>
a <sub>31</sub>	a <sub>32</sub>	a <sub>37</sub>	a <sub>38</sub>	a <sub>33</sub>	a <sub>34</sub>	a <sub>39</sub>	a <sub>35</sub>	a <sub>36</sub>
a <sub>41</sub>	a <sub>42</sub>	a <sub>47</sub>	a <sub>48</sub>	a <sub>43</sub>	a <sub>44</sub>	a <sub>49</sub>	a <sub>45</sub>	a <sub>46</sub>
a <sub>71</sub>	a <sub>72</sub>	a <sub>77</sub>	a <sub>78</sub>	a <sub>73</sub>	a <sub>74</sub>	a <sub>79</sub>	a <sub>75</sub>	a <sub>76</sub>
a <sub>81</sub>	a <sub>82</sub>	a <sub>87</sub>	a <sub>88</sub>	a <sub>83</sub>	a <sub>84</sub>	a <sub>89</sub>	a <sub>85</sub>	a <sub>86</sub>

Local View (CPUs)



# Parallel I/O: MPI I/O

- Use of MPI routines for parallel I/O
- Interfaces for handling shared access on the same file
- Developers must implement a way to store data accordingly to a given data layout => MPI data types helps!



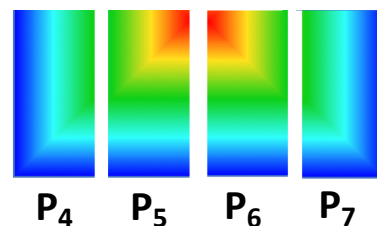
a <sub>11</sub>	a <sub>12</sub>	a <sub>13</sub>	a <sub>14</sub>	a <sub>15</sub>	a <sub>16</sub>	a <sub>17</sub>	a <sub>18</sub>	a <sub>19</sub>
a <sub>21</sub>	a <sub>22</sub>	a <sub>23</sub>	a <sub>24</sub>	a <sub>25</sub>	a <sub>26</sub>	a <sub>27</sub>	a <sub>28</sub>	a <sub>29</sub>
a <sub>31</sub>	a <sub>32</sub>	a <sub>33</sub>	a <sub>34</sub>	a <sub>35</sub>	a <sub>36</sub>	a <sub>37</sub>	a <sub>38</sub>	a <sub>39</sub>
a <sub>41</sub>	a <sub>42</sub>	a <sub>43</sub>	a <sub>44</sub>	a <sub>45</sub>	a <sub>46</sub>	a <sub>47</sub>	a <sub>48</sub>	a <sub>49</sub>
a <sub>51</sub>	a <sub>52</sub>	a <sub>53</sub>	a <sub>54</sub>	a <sub>55</sub>	a <sub>56</sub>	a <sub>57</sub>	a <sub>58</sub>	a <sub>59</sub>
a <sub>61</sub>	a <sub>62</sub>	a <sub>63</sub>	a <sub>64</sub>	a <sub>65</sub>	a <sub>66</sub>	a <sub>67</sub>	a <sub>68</sub>	a <sub>69</sub>
a <sub>71</sub>	a <sub>72</sub>	a <sub>73</sub>	a <sub>74</sub>	a <sub>75</sub>	a <sub>76</sub>	a <sub>77</sub>	a <sub>78</sub>	a <sub>79</sub>
a <sub>81</sub>	a <sub>82</sub>	a <sub>83</sub>	a <sub>84</sub>	a <sub>85</sub>	a <sub>86</sub>	a <sub>87</sub>	a <sub>88</sub>	a <sub>89</sub>
a <sub>91</sub>	a <sub>92</sub>	a <sub>93</sub>	a <sub>94</sub>	a <sub>95</sub>	a <sub>96</sub>	a <sub>97</sub>	a <sub>98</sub>	a <sub>99</sub>

Logical View (Matrix)

a <sub>11</sub>	a <sub>12</sub>	a <sub>17</sub>	a <sub>18</sub>	a <sub>13</sub>	a <sub>14</sub>	a <sub>19</sub>	a <sub>15</sub>	a <sub>16</sub>
a <sub>21</sub>	a <sub>22</sub>	a <sub>27</sub>	a <sub>28</sub>	a <sub>23</sub>	a <sub>24</sub>	a <sub>29</sub>	a <sub>25</sub>	a <sub>26</sub>
a <sub>81</sub>	a <sub>82</sub>	a <sub>87</sub>	a <sub>88</sub>	a <sub>83</sub>	a <sub>84</sub>	a <sub>89</sub>	a <sub>85</sub>	a <sub>86</sub>
a <sub>61</sub>	a <sub>62</sub>	a <sub>67</sub>	a <sub>68</sub>	a <sub>63</sub>	a <sub>64</sub>	a <sub>69</sub>	a <sub>65</sub>	a <sub>66</sub>
a <sub>91</sub>	a <sub>92</sub>	a <sub>97</sub>	a <sub>98</sub>	a <sub>93</sub>	a <sub>94</sub>	a <sub>99</sub>	a <sub>95</sub>	a <sub>96</sub>
a <sub>31</sub>	a <sub>32</sub>	a <sub>37</sub>	a <sub>38</sub>	a <sub>33</sub>	a <sub>34</sub>	a <sub>39</sub>	a <sub>35</sub>	a <sub>36</sub>
a <sub>41</sub>	a <sub>42</sub>	a <sub>47</sub>	a <sub>48</sub>	a <sub>43</sub>	a <sub>44</sub>	a <sub>49</sub>	a <sub>45</sub>	a <sub>46</sub>
a <sub>71</sub>	a <sub>72</sub>	a <sub>77</sub>	a <sub>78</sub>	a <sub>73</sub>	a <sub>74</sub>	a <sub>79</sub>	a <sub>75</sub>	a <sub>76</sub>
a <sub>81</sub>	a <sub>82</sub>	a <sub>87</sub>	a <sub>88</sub>	a <sub>83</sub>	a <sub>84</sub>	a <sub>89</sub>	a <sub>85</sub>	a <sub>86</sub>

Local View (CPUs)

**Good Luck with 3D!!!!**



- You don't need much more of what you have seen yesterday!!
- Add a call to the `H5Pset_dxpl_mpio` to tell HDF5 you are going to use Parallel I/O

```
plist_id = H5Pcreate (H5P_FILE_ACCESS);  
  
hdf5_status = H5Pset_fapl_mpio (plist_id, MPI_COMM_WORLD, MPI_INFO_NULL);
```

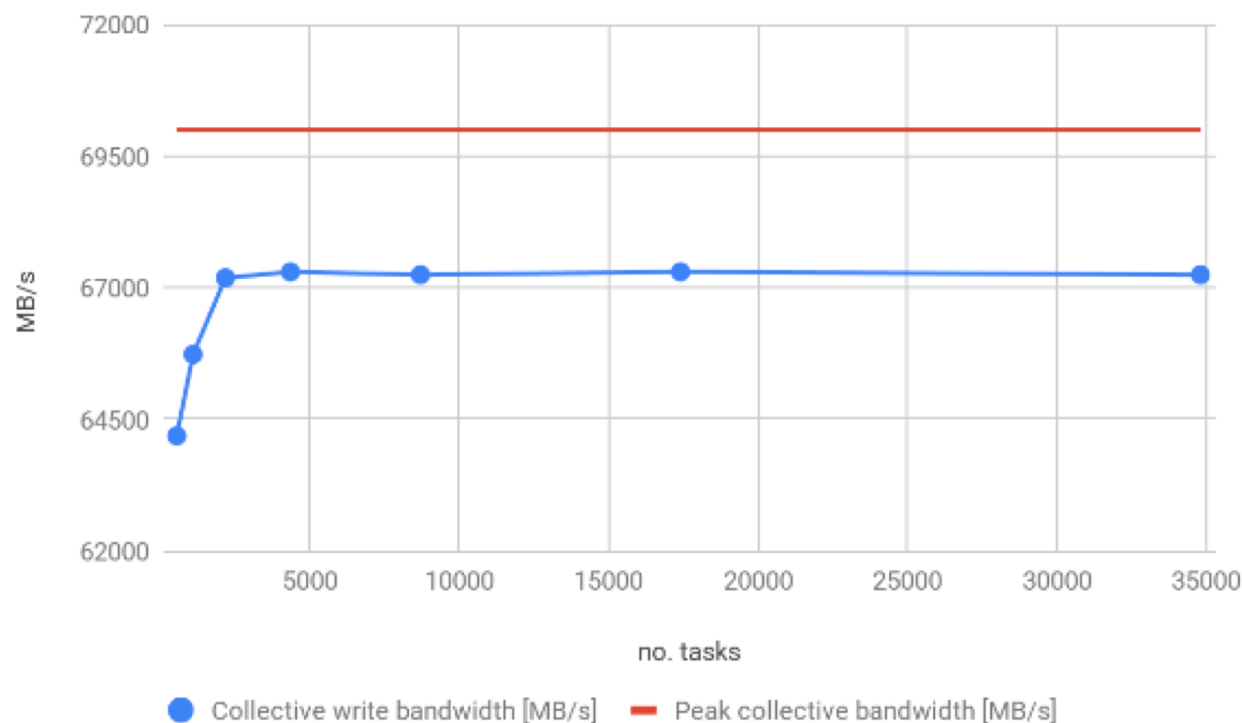
- The hyperslab is still used to describe the distributed dataset
- Remember to define a good balance about what to write in every single file (i.e., trajectories, temporal dynamics, time evolution, ecc. )
  - Few files containing large data
  - A lot of files containing few data

1. Open and close the file every time for writing/reading data (inefficient)
2. Open and close the a different file every time for writing data (inefficient but sometime a good compromise)
3. Open the file for several writing/reading access data
  - ❖ In general more efficient because aiding the I/O system buffering
  - ❖ In writing file can become really big (hard to post-process)
4. Buffer in memory several I/O writing to perform less accesses to the file system with larger dumps from memory to disk (really efficient but requires experience for maintaining a scalable design)

## WHY HDF5

**MAX** DRIVING THE EXASCALE TRANSITION

- API beyond POSIX
- So far the bandwidth usage has shown to be an excellent ratio of the theoretical peak
- HDF5 + MPI-IO is proving to be an efficient, portable and convenient abstraction for high-throughput workloads



\* **Courtesy of Carlo Cavazzoni (CINECA)**