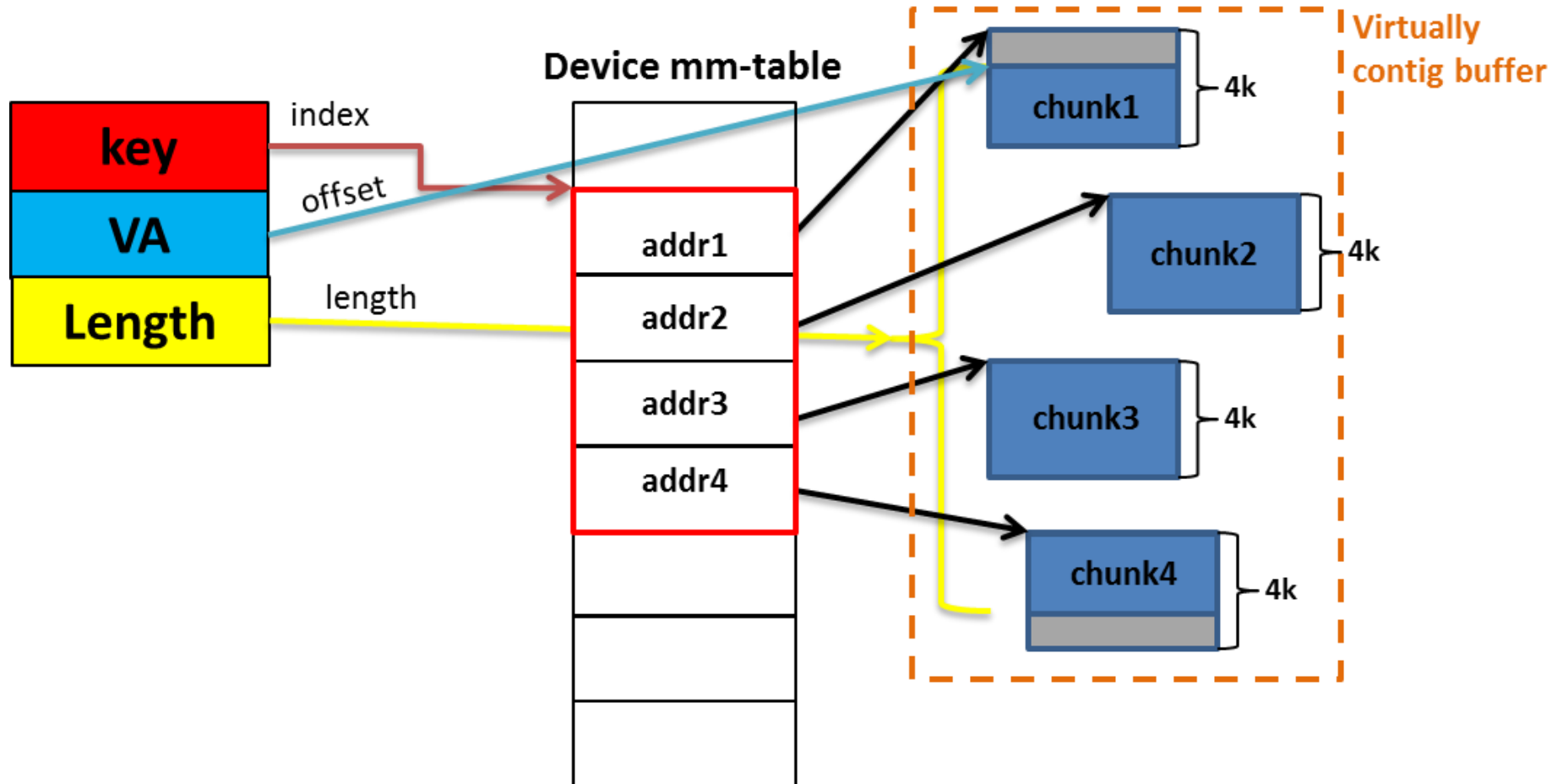




## Indirect Memory registration

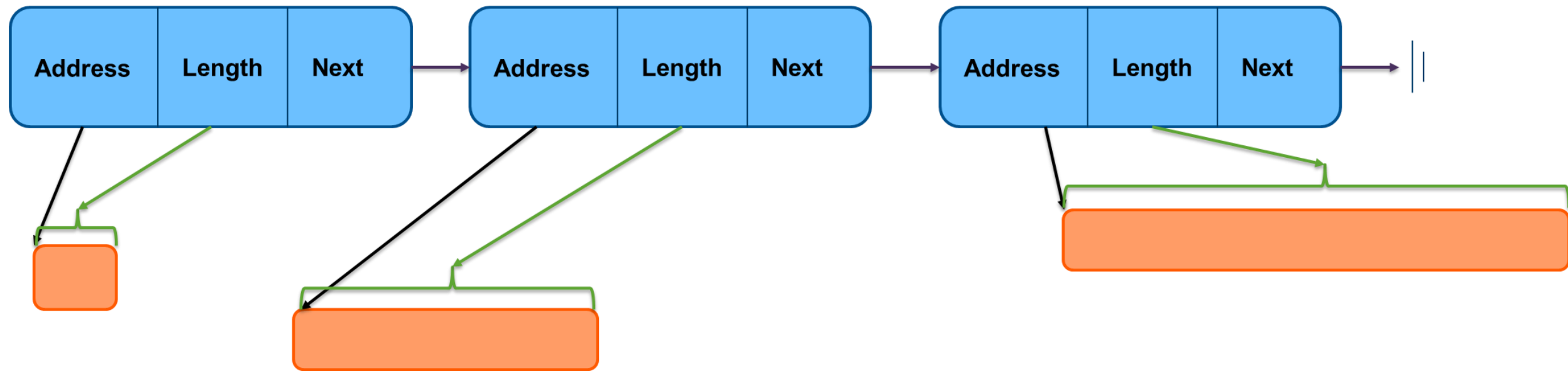
Sagi Grimberg

- RDMA memory registration API allows a specific memory layout (scatter-gather list of buffers)
  - All buffers must be block aligned (single page\_shift)
  - The first buffer is allowed to have a FBO (First Byte Offset)
  - The last buffer is allowed to end before the EOB (End Of Block)





- How can we register an arbitrary SG list of buffers?



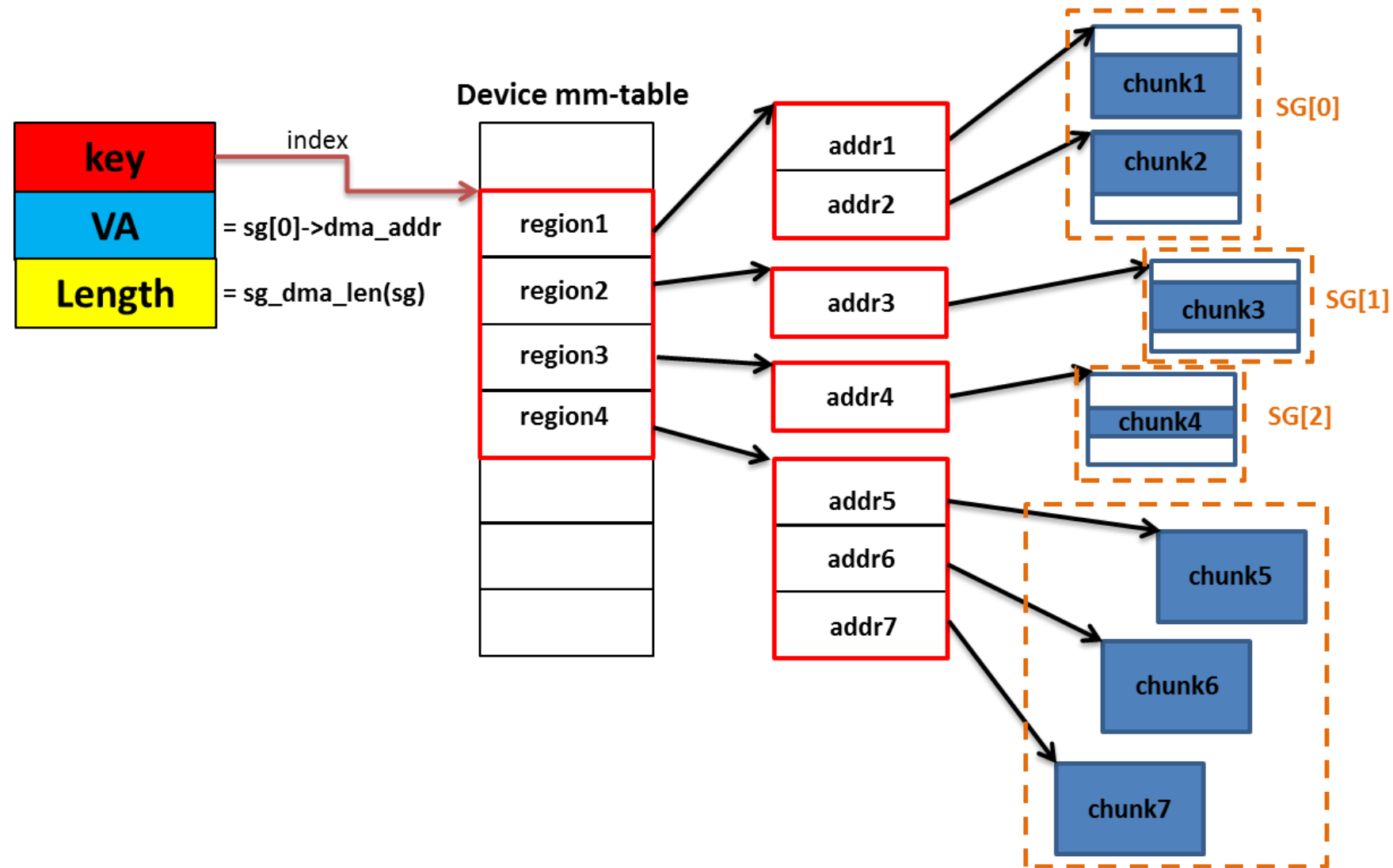
- Possible Solutions:

- Use multiple MRs to register a single SG list (each MR will register an “aligned segment”)
- Use Bounce buffers (SW allocates a well-aligned buffer and copies)
- Break SG list across multiple transactions.
- Use a block size that satisfies the minimal block alignment.

- **Block/file Storage:**
  - Vectored Direct IO
  - SG\_IO
  - Block layer bio merges
  - T10-PI (much easier to get out of alignment in 8-byte DIF buffers merges)
  
- **MPI:**
  - Single RDMA to/from a Remote scatter (in user-space)
  - “User-space FRWR”
  
- **PMEM devices – byte addressable storage (NVM programming model)**
  - We can no longer afford to work-around this limitation and meet the latency requirements.
    - No bounce buffering what-so-ever.
    - Efficient multiple byte range remote access.
  - See [http://www.snia.org/sites/default/files/DougVoigt\\_RDMA\\_Requirements\\_for\\_HA.pdf](http://www.snia.org/sites/default/files/DougVoigt_RDMA_Requirements_for_HA.pdf)
  
- **More?**

# Indirect Memory registration

- Enhance the RDMA stack to allow registering an arbitrary SG lists for devices that support it.



- Follow the same notation as standard FRWR.
  - Instead of passing a page vector to the device we pass a vector of `ib_sge` elements.

OP	FRWR	Indirect
<b>Allocate</b>	<pre>frpl = ib_alloc_fast_reg_page_list() mr = ib_alloc_fast_reg_mr()</pre>	<pre>irl = ib_alloc_indir_reg_list() mr_attrs.create_flags = IB_MR_INDIRECT_REG; mr = ib_create_mr(mr_attrs)</pre>
<b>Free</b>	<pre>ib_free_fast_reg_page_list(frpl) ib_dereg_mr(mr)</pre>	<pre>ib_free_indir_reg_list(irl) ib_dereg_mr(mr)</pre>
<b>Register interface (post_send)</b>	<pre>struct {     u64                iova_start;     struct ib_fast_reg_page_list *page_list;     unsigned int       page_shift;     unsigned int       page_list_len;     u32                length;     int                access_flags;     u32                rkey; } fast_reg;</pre>	<pre>struct {     u64                iova_start;     struct ib_indir_reg_list *indir_list;     unsigned int       indir_list_len;     u64                length;     unsigned int       access_flags;     u32                mkey; } indir_reg;</pre>
<b>Opcode</b>	IB_WR_FAST_REG_MR	IB_WR_REG_INDIR_MR

- Method for generic SG support
  - Synchronous API (e.g. reg\_mr)
  - Asynchronous API (e.g. frwr)
  
- Comments/Feedback on the API itself
  - Combine with FRWR API?
  
- Requirement and place for abstraction





Thank You