# The Case for Energy Clarity

Fan Chung, Henry Kuo, George Candea



### **Energy Usage is Growing Dramatically**



#### **Energy Consumption is not Transparent**



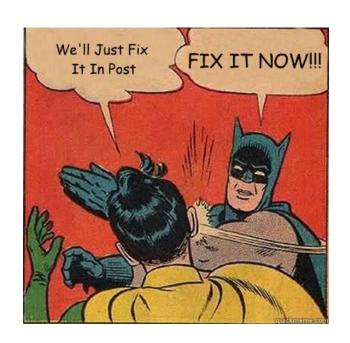
We need energy clarity!

#### **Energy Clarity**

- Ability to reason about energy for all possible inputs:
  - What consumes energy and how much
  - How the consumption changes with the inputs
  - □ ...

#### **Energy Clarity Must Come Before Code**

- Energy clarity shouldn't come as an afterthought.
- Goal: Express, plan, and program energy use **before** the implementation
  - "A call to LZ4\_compress\_default() function should take no more than 35 microjoules, because it's called a lot"
  - "If on the given hardware
     LZ4\_compress\_default() consumes more
     than that energy, then call
     LZ4\_compress\_fast() instead"



But how can we achieve energy clarity?

# **Energy Interfaces!**

#### Having Energy Interfaces Like Having Functional Interfaces

Functional interfaces are introduced in the 1960s to concisely describe a module's functionality

Imagine if we left functional clarity to be an afterthought...

#### Having Energy Interfaces Like Having Functional Interfaces

```
struct list_head *head, **tail = &head; /*
                                            * Returns a list organized in an intermediate
for (;;) {
                                            * format suited to chaining of merge() calls:
     if (cmp(priv, a, b) \leq 0) {
                                            * null-terminated, no reserved or sentinel
         *tail = a;
                                            * head node, "prev" links not maintained.
         tail = &a \rightarrow next;
                                            */
         a = a \rightarrow next:
                                           static struct list head *merge(void *priv,
          if (!a)
                                                list_cmp_func_t cmp, struct list_head *a,
              *tail = b, break;
                                                struct list head *b)
     } else {
         *tail = b;
         tail = \&b \rightarrow next;
         b = b \rightarrow next;
                                         Functional interfaces are key to
          if (!b)
              *tail = a, break;
```

return head;

functional clarity in building systems

**Energy interfaces are key to energy clarity!** 

#### **Energy Interfaces of an ML Web Service**

#### **Energy Interfaces of an ML Web Service**

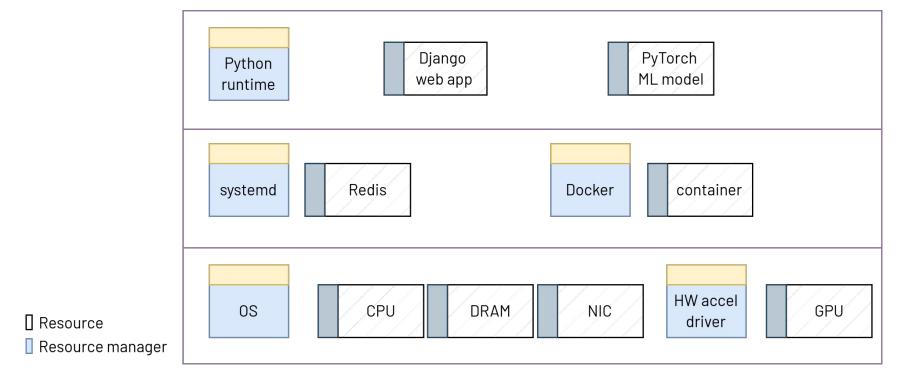
```
def E cache lookup (key, response len):
   # ECV: local cache hit - cache hit in current node
    return 5 if local_cache_hit else 100 * response_len # (Microjoules)
def E cnn forward (image):
    n = 256
    n zeros = image.count(0)
    return (8 * E conv2d(image.size() - n zeros) +
           8 * E relu(n embedding) +
           16 * E mlp(n embedding))
```

#### **Programs as Energy Interfaces**

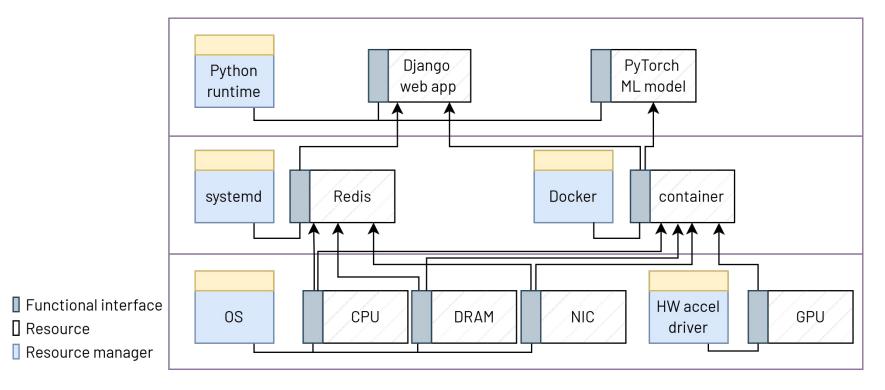
- Energy interfaces are programs that describe how much energy a module will use for any possible input.
  - Appeals to programmers' intuition
  - Executable by programs or tools
  - Precisely describes energy behavior
  - Can express any possible energy behavior

Make energy programmable!

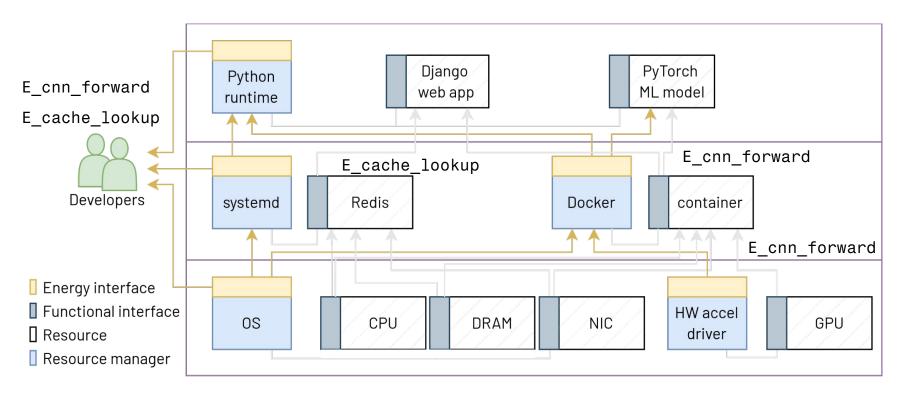
### **Energy Interfaces Exported by Resource mgrs**



### **Energy Interfaces Exported by Resource mgrs**



### **Energy Interfaces Exported by Resource mgrs**



### **Energy Interface** → **Code**

```
Developers
```

```
def E cache lookup (key, response len):
    # ECV: local cache hit - cache hit in current node
    return 5 if local cache hit else 100 * response len # (Microjoules)
const char* cache lookup (const char* key) {
    int i = hash(key);
    if (cache[i].key \&\& strcmp(cache[i].key, key) = 0)
        return cache[i].value;
    const char* val = cache remote fetch(key);
    cache[i] = (CacheEntry){ key, val };
    return val;
```

#### **Code** → **Energy Interface**

```
Changes: Fetching compressed
       const char* val = cache remote fetch(key);
                                                      data and decompressing it locally
       cache[i] = (CacheEntry) { key, val };
       const char* compressed data = cache remote fetch compressed(key);
       const char* val = decompress payload(compressed data);
       cache[i] = (CacheEntry) { key, value };
       return val;
def E_cache lookup(key, response len):
     # ECV: local cache hit - cache hit in current node
     return 5 if local cache hit else 100 * response len
                                                               # (Microjoules)
     return 5 if local cache hit else (70 + 10) * response len # (Microjoules)
```

### **Experiment: GPT-2 Energy Interface**

```
def E GPT 2 inference (wl):
   # ECVs: static power, energy per instr, vram energy per access,
12 energy per access, 11 energy per access
    static_energy = static_power * exec_time[wl.batch sz, wl.precision]
    instruction_energy = energy_per_instr * num_instr[wl.batch_sz, wl.precision]
    vram_energy = vram_energy per_access * sector_rw[wl.batch_sz, wl.precision]
   12_energy = 12_energy_per_access * 12_rw[wl.batch sz, wl.precision]
   11 energy = 11 energy per access * 11 rw[wl.batch sz, wl.precision]
    return static energy + vram energy + 12_energy + 11_energy +
instruction energy
```

### **Preliminary Result**

GPU	Average error	Max error	
RTX4090	0.70%	0.93%	
RTX3070	6.06%	8.11%	

#### Recap

- Energy clarity, as the ability to reason about energy
- Energy interfaces are the key to achieve energy clarity
  - They are **programs** that describe how much energy a module will use for any possible input.
- Resource managers combine and expose the energy interfaces of resources

#### **Open Questions**

#### Modularity

How to divide energy behavior into modules just like functional interfaces?

#### Composability

How to compose the higher-level energy interfaces from lower-level interfaces?

#### Automation

□ Is it possible to automatically obtain the energy interface from a program?

#### Call to Action

- Hardware should export more information about the hardware energy state to software
  - "Energy Counters" equivalent to "Performance Counters"
  - Support for energy measurement with higher frequency, and component-specific granularity.

#### **How Do Energy Interfaces Differ From Other Methods?**

	Energy Profiling	Black-box Energy Modeling	Energy Interface
Can it predict energy consumption?	×		
For any module, can we know how much energy it consumes?			
Can we know how does energy change with inputs?			
Can developers use it to write energy-efficient code?			
Cost-efficiency?			
Feasible to adopt in practice?			Promising)

# **Thank You for Listening!**

# **Thank You for Listening!**