

 Public[Code](#) [Pull requests](#) 312 [Actions](#) [Projects](#) [Security](#) [Insights](#)

master ▾

...

[linux](#) / [arch](#) / [x86](#) / [include](#) / [asm](#) / [msr.h](#)

mdroth x86/boot: Introduce helpers for MSR reads/writes ...

History

18 contributors



384 lines (333 sloc) | 10.8 KB

...

```
1  /* SPDX-License-Identifier: GPL-2.0 */
2  #ifndef _ASM_X86_MSR_H
3  #define _ASM_X86_MSR_H
4
5  #include "msr-index.h"
6
7  #ifndef __ASSEMBLY__
8
9  #include <asm/asm.h>
10 #include <asm/errno.h>
11 #include <asm/cpumask.h>
12 #include <uapi/asm/msr.h>
13 #include <asm/shared/msr.h>
14
15 struct msr_info {
16     u32 msr_no;
17     struct msr reg;
18     struct msr *msrs;
19     int err;
20 };
21
22 struct msr_regs_info {
23     u32 *regs;
24     int err;
25 };
26
27 struct saved_msr {
28     bool valid;
29     struct msr_info info;
```

```
30 };
31
32 struct saved_msrs {
33     unsigned int num;
34     struct saved_msr *array;
35 };
36
37 /*
38  * both i386 and x86_64 returns 64-bit value in edx:eax, but gcc's "A"
39  * constraint has different meanings. For i386, "A" means exactly
40  * edx:eax, while for x86_64 it doesn't mean rdx:rax or edx:eax. Instead,
41  * it means rax *or* rdx.
42  */
43 #ifdef CONFIG_X86_64
44 /* Using 64-bit values saves one instruction clearing the high half of low */
45 #define DECLARE_ARGS(val, low, high)    unsigned long low, high
46 #define EAX_EDX_VAL(val, low, high)    ((low) | (high) << 32)
47 #define EAX_EDX_RET(val, low, high)    "a" (low), "d" (high)
48 #else
49 #define DECLARE_ARGS(val, low, high)    unsigned long long val
50 #define EAX_EDX_VAL(val, low, high)    (val)
51 #define EAX_EDX_RET(val, low, high)    "A" (val)
52 #endif
53
54 /*
55  * Be very careful with includes. This header is prone to include loops.
56  */
57 #include <asm/atomic.h>
58 #include <linux/tracepoint-defs.h>
59
60 #ifdef CONFIG_TRACEPOINTS
61 DECLARE_TRACEPOINT(read_msr);
62 DECLARE_TRACEPOINT(write_msr);
63 DECLARE_TRACEPOINT(rdpmc);
64 extern void do_trace_write_msr(unsigned int msr, u64 val, int failed);
65 extern void do_trace_read_msr(unsigned int msr, u64 val, int failed);
66 extern void do_trace_rdpmc(unsigned int msr, u64 val, int failed);
67 #else
68 static inline void do_trace_write_msr(unsigned int msr, u64 val, int failed) {}
69 static inline void do_trace_read_msr(unsigned int msr, u64 val, int failed) {}
70 static inline void do_trace_rdpmc(unsigned int msr, u64 val, int failed) {}
71 #endif
72
73 /*
74  * __rdmsr() and __wrmsr() are the two primitives which are the bare minimum MSR
75  * accessors and should not have any tracing or other functionality piggybacking
76  * on them - those are *purely* for accessing MSRs and nothing more. So don't even
77  * think of extending them - you will be slapped with a stinking trout or a frozen
78  * shark will reach you, wherever you are! You've been warned.
```

```
79  */
80  static __always_inline unsigned long long __rdmsr(unsigned int msr)
81  {
82      DECLARE_ARGS(val, low, high);
83
84      asm volatile("1: rdmsr\n"
85                  "2:\n"
86                  _ASM_EXTABLE_TYPE(1b, 2b, EX_TYPE_RDMSR)
87                  : EAX_EDX_RET(val, low, high) : "c" (msr));
88
89      return EAX_EDX_VAL(val, low, high);
90  }
91
92  static __always_inline void __wrmsr(unsigned int msr, u32 low, u32 high)
93  {
94      asm volatile("1: wrmsr\n"
95                  "2:\n"
96                  _ASM_EXTABLE_TYPE(1b, 2b, EX_TYPE_WRMSR)
97                  : : "c" (msr), "a"(low), "d" (high) : "memory");
98  }
99
100 #define native_rdmsr(msr, val1, val2) \
101 do { \
102     u64 __val = __rdmsr((msr)); \
103     (void)((val1) = (u32)__val); \
104     (void)((val2) = (u32)(__val >> 32)); \
105 } while (0)
106
107 #define native_wrmsr(msr, low, high) \
108     __wrmsr(msr, low, high)
109
110 #define native_wrmsrl(msr, val) \
111     __wrmsr((msr), (u32)((u64)(val)), \
112             (u32)((u64)(val) >> 32))
113
114 static inline unsigned long long native_read_msr(unsigned int msr)
115 {
116     unsigned long long val;
117
118     val = __rdmsr(msr);
119
120     if (tracepoint_enabled(read_msr))
121         do_trace_read_msr(msr, val, 0);
122
123     return val;
124 }
125
126 static inline unsigned long long native_read_msr_safe(unsigned int msr,
127                                                         int *err)
```

```

128 {
129     DECLARE_ARGS(val, low, high);
130
131     asm volatile("1: rdmsr ; xor %[err],%[err]\n"
132                 "2:\n\t"
133                 _ASM_EXTABLE_TYPE_REG(1b, 2b, EX_TYPE_RDMSR_SAFE, %[err])
134                 : [err] "=r" (*err), EAX_EDX_RET(val, low, high)
135                 : "c" (msr));
136     if (tracepoint_enabled(read_msr))
137         do_trace_read_msr(msr, EAX_EDX_VAL(val, low, high), *err);
138     return EAX_EDX_VAL(val, low, high);
139 }
140
141 /* Can be uninline because referenced by paravirt */
142 static inline void notrace
143 native_write_msr(unsigned int msr, u32 low, u32 high)
144 {
145     __wrmsr(msr, low, high);
146
147     if (tracepoint_enabled(write_msr))
148         do_trace_write_msr(msr, ((u64)high << 32 | low), 0);
149 }
150
151 /* Can be uninline because referenced by paravirt */
152 static inline int notrace
153 native_write_msr_safe(unsigned int msr, u32 low, u32 high)
154 {
155     int err;
156
157     asm volatile("1: wrmsr ; xor %[err],%[err]\n"
158                 "2:\n\t"
159                 _ASM_EXTABLE_TYPE_REG(1b, 2b, EX_TYPE_WRMSR_SAFE, %[err])
160                 : [err] "=a" (err)
161                 : "c" (msr), "0" (low), "d" (high)
162                 : "memory");
163     if (tracepoint_enabled(write_msr))
164         do_trace_write_msr(msr, ((u64)high << 32 | low), err);
165     return err;
166 }
167
168 extern int rdmsr_safe_regs(u32 regs[8]);
169 extern int wrmsr_safe_regs(u32 regs[8]);
170
171 /**
172  * rdtsc() - returns the current TSC without ordering constraints
173  *
174  * rdtsc() returns the result of RDTSC as a 64-bit integer. The
175  * only ordering constraint it supplies is the ordering implied by
176  * "asm volatile": it will put the RDTSC in the place you expect. The

```

```
177  * CPU can and will speculatively execute that RDTSC, though, so the
178  * results can be non-monotonic if compared on different CPUs.
179  */
180  static __always_inline unsigned long long rdtsc(void)
181  {
182      DECLARE_ARGS(val, low, high);
183
184      asm volatile("rdtsc" : EAX_EDX_RET(val, low, high));
185
186      return EAX_EDX_VAL(val, low, high);
187  }
188
189  /**
190   * rdtsc_ordered() - read the current TSC in program order
191   *
192   * rdtsc_ordered() returns the result of RDTSC as a 64-bit integer.
193   * It is ordered like a load to a global in-memory counter. It should
194   * be impossible to observe non-monotonic rdtsc_unordered() behavior
195   * across multiple CPUs as long as the TSC is synced.
196   */
197  static __always_inline unsigned long long rdtsc_ordered(void)
198  {
199      DECLARE_ARGS(val, low, high);
200
201      /*
202       * The RDTSC instruction is not ordered relative to memory
203       * access. The Intel SDM and the AMD APM are both vague on this
204       * point, but empirically an RDTSC instruction can be
205       * speculatively executed before prior loads. An RDTSC
206       * immediately after an appropriate barrier appears to be
207       * ordered as a normal load, that is, it provides the same
208       * ordering guarantees as reading from a global memory location
209       * that some other imaginary CPU is updating continuously with a
210       * time stamp.
211       *
212       * Thus, use the preferred barrier on the respective CPU, aiming for
213       * RDTSCP as the default.
214       */
215      asm volatile(ALTERNATIVE_2("rdtsc",
216                                "lfence; rdtsc", X86_FEATURE_LFENCE_RDTSC,
217                                "rdtscp", X86_FEATURE_RDTSCP)
218                  : EAX_EDX_RET(val, low, high)
219                  /* RDTSCP clobbers ECX with MSR_TSC_AUX. */
220                  :: "ecx");
221
222      return EAX_EDX_VAL(val, low, high);
223  }
224
225  static inline unsigned long long native_read_pmc(int counter)
```

```

226 {
227     DECLARE_ARGS(val, low, high);
228
229     asm volatile("rdpmc" : EAX_EDX_RET(val, low, high) : "c" (counter));
230     if (tracepoint_enabled(rdpmc))
231         do_trace_rdpmc(counter, EAX_EDX_VAL(val, low, high), 0);
232     return EAX_EDX_VAL(val, low, high);
233 }
234
235 #ifdef CONFIG_PARAVIRT_X86
236 #include <asm/paravirt.h>
237 #else
238 #include <linux/errno.h>
239 /*
240  * Access to machine-specific registers (available on 586 and better only)
241  * Note: the rd* operations modify the parameters directly (without using
242  * pointer indirection), this allows gcc to optimize better
243  */
244
245 #define rdmsr(msr, low, high) \
246 do { \
247     u64 __val = native_read_msr((msr)); \
248     (void)((low) = (u32)__val); \
249     (void)((high) = (u32)(__val >> 32)); \
250 } while (0)
251
252 static inline void wrmsr(unsigned int msr, u32 low, u32 high)
253 {
254     native_write_msr(msr, low, high);
255 }
256
257 #define rdmsrl(msr, val) \
258 ((val) = native_read_msr((msr)))
259
260 static inline void wrmsrl(unsigned int msr, u64 val)
261 {
262     native_write_msr(msr, (u32)(val & 0xffffffffULL), (u32)(val >> 32));
263 }
264
265 /* wrmsr with exception handling */
266 static inline int wrmsr_safe(unsigned int msr, u32 low, u32 high)
267 {
268     return native_write_msr_safe(msr, low, high);
269 }
270
271 /* rdmsr with exception handling */
272 #define rdmsr_safe(msr, low, high) \
273 ({ \
274     int __err; \

```

```

275     u64 __val = native_read_msr_safe((msr), &__err);          \
276     (*low) = (u32)__val;                                       \
277     (*high) = (u32)(__val >> 32);                             \
278     __err;                                                     \
279 })
280
281 static inline int rdmsrl_safe(unsigned int msr, unsigned long long *p)
282 {
283     int err;
284
285     *p = native_read_msr_safe(msr, &err);
286     return err;
287 }
288
289 #define rdpmc(counter, low, high)                               \
290 do {                                                            \
291     u64 _l = native_read_pmc((counter));                       \
292     (low) = (u32)_l;                                           \
293     (high) = (u32)(_l >> 32);                                  \
294 } while (0)
295
296 #define rdpmcl(counter, val) ((val) = native_read_pmc(counter))
297
298 #endif /* !CONFIG_PARAVIRT_XXL */
299
300 /*
301  * 64-bit version of wrmsr_safe():
302  */
303 static inline int wrmsrl_safe(u32 msr, u64 val)
304 {
305     return wrmsr_safe(msr, (u32)val, (u32)(val >> 32));
306 }
307
308 struct msr *msrs_alloc(void);
309 void msrs_free(struct msr *msrs);
310 int msr_set_bit(u32 msr, u8 bit);
311 int msr_clear_bit(u32 msr, u8 bit);
312
313 #ifdef CONFIG_SMP
314 int rdmsr_on_cpu(unsigned int cpu, u32 msr_no, u32 *l, u32 *h);
315 int wrmsr_on_cpu(unsigned int cpu, u32 msr_no, u32 l, u32 h);
316 int rdmsrl_on_cpu(unsigned int cpu, u32 msr_no, u64 *q);
317 int wrmsrl_on_cpu(unsigned int cpu, u32 msr_no, u64 q);
318 void rdmsr_on_cpus(const struct cpumask *mask, u32 msr_no, struct msr *msrs);
319 void wrmsr_on_cpus(const struct cpumask *mask, u32 msr_no, struct msr *msrs);
320 int rdmsr_safe_on_cpu(unsigned int cpu, u32 msr_no, u32 *l, u32 *h);
321 int wrmsr_safe_on_cpu(unsigned int cpu, u32 msr_no, u32 l, u32 h);
322 int rdmsrl_safe_on_cpu(unsigned int cpu, u32 msr_no, u64 *q);
323 int wrmsrl_safe_on_cpu(unsigned int cpu, u32 msr_no, u64 q);

```

```
324 int rdmsr_safe_regs_on_cpu(unsigned int cpu, u32 regs[8]);
325 int wrmsr_safe_regs_on_cpu(unsigned int cpu, u32 regs[8]);
326 #else /* CONFIG_SMP */
327 static inline int rdmsr_on_cpu(unsigned int cpu, u32 msr_no, u32 *l, u32 *h)
328 {
329     rdmsr(msr_no, *l, *h);
330     return 0;
331 }
332 static inline int wrmsr_on_cpu(unsigned int cpu, u32 msr_no, u32 l, u32 h)
333 {
334     wrmsr(msr_no, l, h);
335     return 0;
336 }
337 static inline int rdmsrl_on_cpu(unsigned int cpu, u32 msr_no, u64 *q)
338 {
339     rdmsrl(msr_no, *q);
340     return 0;
341 }
342 static inline int wrmsrl_on_cpu(unsigned int cpu, u32 msr_no, u64 q)
343 {
344     wrmsrl(msr_no, q);
345     return 0;
346 }
347 static inline void rdmsr_on_cpus(const struct cpumask *m, u32 msr_no,
348                                 struct msr *msrs)
349 {
350     rdmsr_on_cpu(0, msr_no, &(msrs[0].l), &(msrs[0].h));
351 }
352 static inline void wrmsr_on_cpus(const struct cpumask *m, u32 msr_no,
353                                 struct msr *msrs)
354 {
355     wrmsr_on_cpu(0, msr_no, msrs[0].l, msrs[0].h);
356 }
357 static inline int rdmsr_safe_on_cpu(unsigned int cpu, u32 msr_no,
358                                     u32 *l, u32 *h)
359 {
360     return rdmsr_safe(msr_no, l, h);
361 }
362 static inline int wrmsr_safe_on_cpu(unsigned int cpu, u32 msr_no, u32 l, u32 h)
363 {
364     return wrmsr_safe(msr_no, l, h);
365 }
366 static inline int rdmsrl_safe_on_cpu(unsigned int cpu, u32 msr_no, u64 *q)
367 {
368     return rdmsrl_safe(msr_no, q);
369 }
370 static inline int wrmsrl_safe_on_cpu(unsigned int cpu, u32 msr_no, u64 q)
371 {
372     return wrmsrl_safe(msr_no, q);
```



```
373     }
374     static inline int rdmsr_safe_regs_on_cpu(unsigned int cpu, u32 regs[8])
375     {
376         return rdmsr_safe_regs(regs);
377     }
378     static inline int wrmsr_safe_regs_on_cpu(unsigned int cpu, u32 regs[8])
379     {
380         return wrmsr_safe_regs(regs);
381     }
382     #endif /* CONFIG_SMP */
383     #endif /* __ASSEMBLY__ */
384     #endif /* _ASM_X86_MSR_H */
```