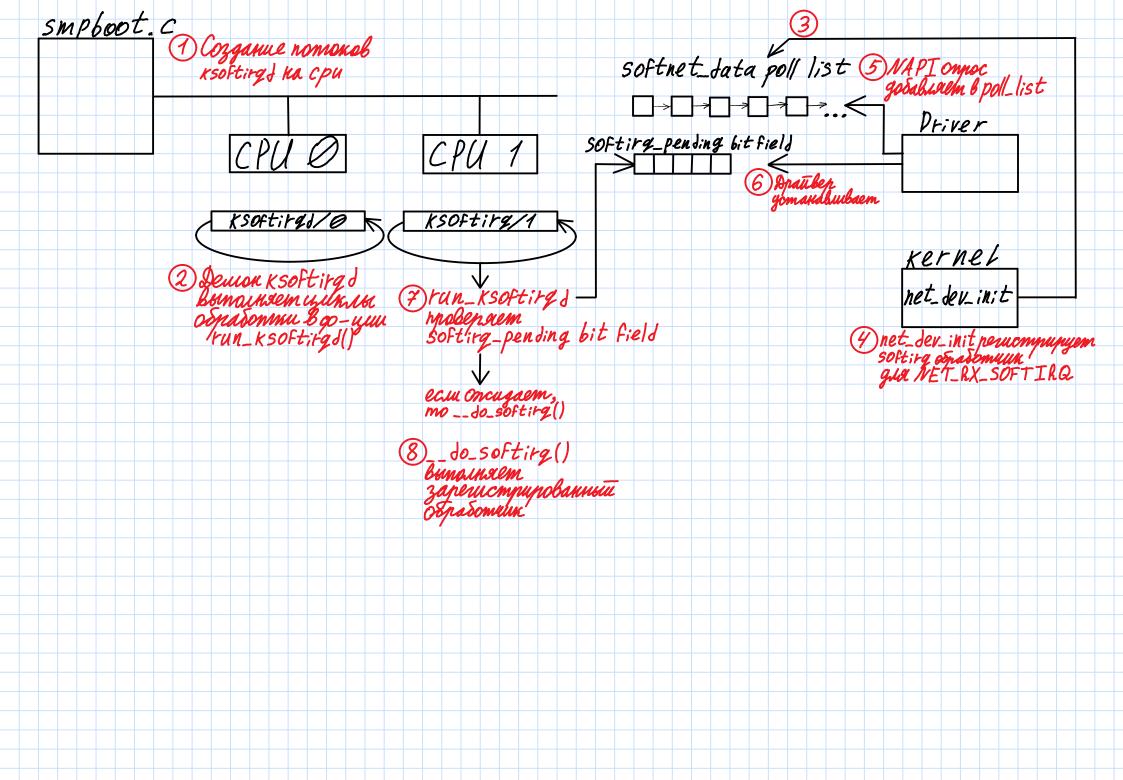
Прозапсение про аппаранные прерывания (Tyungun pacnapamembanua opynnyui) B colp. OC noabiniacs cemebase nogcuemena, neu (man monce econs yenn-ba a out monce sibi. внешний)и отправичить информациюм наulso kalmbromera; Обрабонска прер-и-ванснеймая задача
мобей ОС (по крайней мере, общего низка-Menual; Commence amapaners men-u garacio Commensione na Cocconer ynobran mugu-Сиопсиость обрадонни таких пред-и с текошими времени уванчиванась, что привено к уваничению времени выполнения обработников прер-и Прерватовот-е такого обработима негода, on bomainseman om rayana go nouga; Это очень сильно вличет на отзывнивость системы; Sanbuol breus — breus o scryncubanus zampo-ca breu. ycmp-ba/morsecca cmporo ompegareno — — ncecmuol peanonol breus;

Omzoblubecons cucmenta + peansone brenne; (peanyur na genembur) (+ Enempo) (reakyur na genembur) au. onp-e pearsono brenenu POSIX;} Omzobulbocomo cucmenios pregnonaralm branco Ombema na genembina norozobamena, nomopoe ue mubaguno do n encuganuo u mo ne reasonel enema; Ображиний прер-и домпены завершаться Kak Moncho Esiempee, b OC marie ospasomunu men-u conam genumo na 2 naisbuna: top half ubottom half; В рез- не обработник прер-я дажен выhax mabuno, mo nouqueme gambo (ecu mo bbog) uz <u>будоера контранера устр-ва</u> и eë передача в будоер ядра; Marion osposoniur men-a repeg chowe zabepwenuen garnen unuguaruzunobams omincennoe generabue; Omyonceunice genembur b colop acomenax Subarom 3 bugob: -softirg; -tasklet; -workqueue;

softirg omegeneum communecus; 10 oбработников, представлены в массиве: const char * softifg_to-name[NR_SOFT_IRQS]= {"HI", "TIMER", "NET_Tx", "NET_kx", "BLOCK",
"BLOCK_IOPOLL", "TASKLET", "SCHED", "HR_TIMER", " KCU" 3; Pla softirg ameglulus 90-yuu agpa u cucm. berzeber; struct softing_action void (*action)(struct softing_action*); Korga signo bomounsem oбрабонщик omnoncennoro прерывания, cp-yus act;on() buzubalmis c yrazamenen na struct soft;rq_action b rarecmbe арпичента; extern void open_softirg(intur,
void (*action(struct soft;rg_action **)); ma go-yun unuyuaunyunyem softirg; nr-ungeke al-ma maccuba softity_vec; action-grazaments na go-yuno softirg, Komopan Sygem Comonnamocn; softirg_to_vec onpegeien anauoununo softifg_to-namé: mo maccub pagniepem 10 21-mob, l' romonau repenuement une a copa-Somuncos, Haruhan c'e;

Ospasonnun omnoncennon genembur (softirg), nom-poe zapernompupobano 90-yuen open_softirg(), b poz-me bun-r mon 90-yun comabumen b'orenego na bomornemie; mo on concernoe genembre annubusupyenca c nousely 600 go-you raise_softing: extern void raise_softing (unsigned int nr) local_irg_restore(flags); это макросы; 9-yun raise_softirg_irqoff() мопсет выпол-няться только с разремёнными терываниями; Strumman gin nancgow mouseccopa cynsecm-byem k softirg Jaemon-gluon, patomarousia c omioncennum gluomburum (softirg); Ha примере получения сетевых пакетов проземонстрируем посл-пь денствый: Truzanycke cucmeuk gik kancgoro egra hanukalm Benraumung gewon softite d; ! 3 Imo yuni npobernu; бакие-то вкешние действия (в настности, инициализация девайса, связанного с сетью) тиводит к формированию события в системе (в данной случае это получение сетевых паке-тов), которое дамско быть обработамо; Coomb. npep-e, komopol npugiem om cemeboro ycmp-ba, rpublegiem k zamamenino m.n. po///ist u zameni c naucusho softirg daemon, na komopon bygem bunomena go-una run_ksoftirgd, sma ungoopma-una unusuamusupulm busob go-unu — do_softirg();



B repeaul nomme irq qua cemebou nogenements upucymembyem maunep (HI). Beë mo bunamaemas report softing daemon u burgob coomb. go-yun Agna. Truxogaujue annapamuice npep-a innugua-Mizuryrom oncuganul bemonnenia zaperucnimpobartion omioncennos genembus; NAPI - New API (nouberrocs & bencu agra 2.4) Toaburca, buguno, kan ombem na bozpocuniu orben genembun c cemebon nogenementon; Они разработан для улучиения выполнения (да, умучиения, это цитата) высокоскоростных gencinbun c cembon rogencmenson; (Mas. High speed networking) High speed networking woncern renepupolami (согдавать) тисячи препиваний в секупру, и они домски быть обработаны; Очевидно, что они обрабатываются именно таким образом: спачана возмикшее прер-е peruconjungemen, zamen navienzalmen B coomb. orenego u obnavamenbalemen l'ecomb-un e mou orenegon; (xypcoban?) {Cemebaa nogeucmena abn. kanbonenemen }
gia cepbepeb;

101- momokai naginonoso coegunemur, т.к. троиного рукопопсания: чтобы начать opavambani nakem, kulem zanjanubalm cepben, cepben nocuiaem ombemuse coorigenue u negëm nogmbepnegenur om muluma. Konga nogmblepricaline om kulenna ke njuxogum, соединение наз. полутирымым. Но такое coequienue zasubalm bxoguyo orepegs, u 6 rezyromane remmunice navemin ne ospadamillaromaa; (DDOS - amanu) raise_softirg_irg_off() naueuaem softirg kan omnoncennu nymën yemanoban coomb. Tuma b Sumobou uache softing (softing pending novambuoso npoyec-copa). Imogenalman b go-yun raise_softing (); East you man mosseccop bunawaem annapamioe nnen-e, mo ocyineconsulueman buxog uz go-yuu raise_softirg_irgoff() u boccmanabмвается дагах IF. в топпивной смучае (если не в прерывании) вызывается wakeup_softirg d():

if (!in_interrupt)
waxeup_softirgd();

Imo: Swacmen, nouling kancestu undereccop undem coscmbenniu Ksoftirg Saemon;

Kancgun softing monogum regres areg. mansi: 1) softirg pupyemen go-yueu Open_softirg(); 2) Oppationium softing omwellalman kar omnonclusion (nousous to go-your raise-softing); 3 Bamen ommelleunie soft it g zangenatomas na bomainence (triggered) buancjou areginusen yume bonownehun omnoncentura 4) Bunonience zanamubalmar zabepulencen; Уштата из источника (макуал!): "Toncaujuema, uzbecatume borgerenun nobux softirg's, ecuu bam ne nyneno reasono borconocuonocmuse nomonobol manipolanul zagamin. L'é b quiunail, bisconoracmonne, m. é. écui Re HYNCHO C BUCONOU L'ACMONION BUMANNO manupobanue omnoncemmen genembin; Dun normu beex nynenoux yenen dance nem germaniouno tasklet'ob (onu onnegeneum 6 10-mu soft :rg); tasklet- cnew mun soft itq; Hannunen, bee bot tom half nocuegobamensmux yomnouemb gamenn Somb reasuzobann nan tasklet'u, a ne nan softirg; tasklet-unoronomornom anavor bot tom half;

Удно версии 6.2.2: "Invom API yomapen. Toncavyucma, paccinom-pume bozuonenocorb uenavozobanua nomono-bux irg-zanpocob buecmo muoro." Umo za nomonobole i rg-zanpocu! Chopel belio peut ugėm oб очередях работ (workqueue); & ommune om softirg, tasklet woncem on normous montoso ra ognou monseccone, M.l. gur macmiensob regonyemmino ux napan-Missiol bomoinence (pero ugemo napamentoman bomoinence ognorou moso nee macmiena); Combenancembenno, ir q morym bornounsunca napanrenono (au. Kapmunky); Для тастения установнено отаншение: Ogun u nom nel macmiem moncem bonnomoca такого в единственным экзентияре; A softirg gonycharom napamelubriol bomonnemie => ybenivence cropsonin odaynculanun manux zanpocos; M.K. ogun u mom nel mun softirg moncem bemon-RAMICA (planono) napamentoso, mo na nux ramagos-baronca ncicmane mperobarua no cicnonozobaruso chedomp prannaicum monthus;

II. K. ogun u mom nee tasklet ne moncem Bomonnombca napamentono, mo manux ncecmmux mperobania na tasklet'u ne namagosbalmaa; If m mon tasklet he woncen aconyposamica i bomanaeman om navara go noniga; Ila tasklet'ax amegenena compyunyma (v.6.2.2) struct tasklet_struct //ecomb Buemognuke struct tasklet_struct *next; //concentaskletics unsigned long state; atomic-t count; bool use-callback; Bräupozuwa? void (*func)(unsigned long data);

void (*callback)(struct tasklet_struct *t); unsigned long data; tasklet'u monym penicmpupobamoca b cucmene KAK CMAMWECKU, MAK U GUHANWECKU; Для статического объявления тактета суuzecmbyem 2 Manpoca: DECLARE_TASKLET(name,_callback); DECLARE_TASKLET_DISABLED(name,_callback); Cyuzecmbyrom mannce obsorbienua m.n. cmapus tasklet'ob. B nux bruecmo callback ucnalozyemca

```
func:
DECLARE_TASKLET_OLD (name, _ func);
DECLARE_TASKLET_DISABLED_OLD (name, _ func);
    Dunanune obsesseence toesklet a:
extern void tasklet_iuit (struct tasklet_struct *t, void (*func)(unsigned long), unsigned long tata);
   Tare moro, kan tæsklet obuselnen, on gameln bumb
zamanupoban na bunamenue:
static inline voit tasklet_schedule(
struct tasklet_struct **t)
       if (!test_and_set_bit(TASKLET_STATE_SCHED,
t->state))
                   tasklet_schedule();
Static inline tasklet_hi_schedule(...);

Suconomnommemme tasklet'u sygym

bomomumban panome oswnox;
   B Saree rannux cucmentax econs euze tasklet
schedule_f;rst, no b colp. Agre mon qo-yun nem;
   D-yuno tasklet_schedule() nynono borzbanib
bornacomune npepubanua;
```

Cb-ba tasklet 106: 1/ Eau buzbana op-yun tasklet_schefule(), mo nocue eë buzoba tasklet ranannynobanno dygem bemoissell karau-moo mousecope xome on ogun paz, 2) Ecu tasklet zamanupoban, no enzë ne bomonnaeconcer, on Sygem bomounted maious ogun paz; 3) Ecu tasklet ynce bonainsemen na gpijnou npoyeccope ulli go-yun tasklet_schefule() Bizoi-Baemaa uz causio tasket'a, mo Bomomenne oyaem omnonceno na carel nozgnun crox; 4) Eau pagpaonyux cumaen, uno b tasxletie resorrogumo bomonumo genembus, comopore mannee monym bomanniben b gpyrom tasklet'e, mo b tasklet'e kago ucnollzolanio bzalluouckinorence c nauouyoto spin lock; MORCHEMUE: tasklet bemainsemen om navara go nomya morbro na ogkow npolseccope (napawellbro bunormounted He moncem, normouly que nouncem uno tasklet'a bzammonomnovenice не нупско. Lo eau parnasomeur negnonoraem, umo gba poznoux tasklet'a M. d. zaummenecobanu B ognux u mex nel garmoux, mob nogax enux tasklet'ob датско быть реализовано взаимоистиочение c nouseurs spin lock; 13 Mo, umo tasklet'y ne mpesylmen numanux cpegemb byan-uoucannovenus no omn-ro k canony cese, b grune-

13 nave nazorbalnica cepuaruzarzulu; Ha taskletian omnegeneru go-yun tasklet_disable() u tasklet_enable(); struct tasklet_struct B colp. nogax omnimuence tasklet ob uz amouan-HUX (Regenmon) KOMMENCMOB Repelano oumo-Kanu, u этого следует избегать; tasklet- x: 11(struct tasklet_struct *); tasxlet_setup(); BI/p: tasklet manunyemen bospasom une amanammon mensibanun; Tianupobanue tasklet'a moncem bomonnoman 6 90-yuu mogym module_init(); Обработник аппаратного прер-я надо реinempupobamio na pazgeninemoti nimuira; Dias IRQF_SHARED nozbourem zaperucnyu pobamil nam ornadomnim npen-a na minum itq, chazannoti c maduamyroti; Cettrac IhQ12 re bugno, no bugno IRQ8(API); Pansue cyuzecomboban oney gaiar Europoro nplp-a(quick). Courac equinemb. Sucmpun prep-en comanace npep-e om cucm. mannepa;

Ettpexuunnaa aprumennypa? Honazamo men-a moncho repez /proc/interrupts Resocraguiso resecums o Epasonilik men-a na mumio Ih Q1 (9-e men-e); Sleedrogum namicamb odpadomina annapam-noro rpep-a, komopini dygem manupobami bu-namenie omnoncemoro generalia b buge tasklet'a; task let gonneen bubogums breun, morga On Sun zamanypoban u konga sun bomannen; Fax mablino, tasklet bomonuseman na mon псе процессоре, на которым выполнямся об-работник прер-я, которым его запланирован; в наших архименнурах, когда выполняется аппаратное прер-е, на том ироцессоре, на ко-тором выполняется обработник прер-я, все nnen-a zanpouseusi; Ab cucmene zanpenzena men-a na bcex npo-neccopax no gannon summe Ih a;