# Construction of an EST database for *Bombyx mori* and its applications

**Kazuei Mita**[*,†], **Mitsuoki Morimyo**[‡], **Kazuhiro Okano**[§], **Yoshiko Koike**[#], **Junko Nohata**[†], **Masataka G. Suzuki**[#] **and Toru Shimada**[#]

[†]Laboratory of Insect Genome, National Institute of Agrobiological Sciences, Owashi 1-2, Tsukuba, Ibaraki 305-8634, Japan
[‡]National Institute of Radiological Sciences, Anagawa 4-9-1, Inage-ku, Chiba 263-8555, Japan
[§]Laboratory of Molecular Entomology and Baculovirology, The Institute of Physical and Chemical Research (RIKEN), Hirosawa 2-1, Wako, Saitama 351-0198, Japan
[#]Department of Agricultural and Environmental Biology, University of Tokyo, Yayoi 1-1-1, Bunkyo-ku, Tokyo 113-8657, Japan

To build a foundation for the genome analysis of *Bombyx mori*, we have constructed an EST database. Since gene expression patterns deeply depend on tissues as well as developmental stages, we analysed many cDNA libraries prepared from various tissues and different developmental stages to cover the whole set of *Bombyx* genes. So far, the *Bombyx* EST db contains 26,000 ESTs from 32 cDNA libraries, which are grouped into about 9,500 non-redundant ESTs. Based on the EST db, we prepared a microarray containing 6,000 non-redundant ESTs for functional studies as well as for genome analysis. In addition, we have constructed a high quality BAC library for *Bombyx mori* characterized by an average insertion size of about 170 kb with 11-fold redundancy. We are constructing *Bombyx* BAC contigs by filter hybridization using non-redundant ESTs as probes.

IN eukaryotes, the genome projects of various species have been vigorously pushed forward; genome sequences have been completed in yeast, *Caenorhabditis elegans*, *Arabidopsis* and *Drosophila*, while those of human and mouse will be finished shortly. Genoscope and Celera are carrying out *Anopheles* genome analysis. In addition, whole genome sequencing is planned in many other species. Genome information provides powerful tools to understand biological mechanisms and functions and is essential for biology, medical science and agriculture.

The Lepidoptera to which the silkworm belongs, include the most highly destructive agricultural pests; hundreds of species of caterpillars cause widespread economic damage on food and fibre crop plants, fruit trees, forests and stored grains. Lepidopteran genome information must make a strong impact on insect science and industries such as insecticide, pest control, and silk production (see International Lepidopteran Genome Project Proposal, 2001)[1]. In Lepidoptera, however, genome information is quite limited so far. As *Drosophila* is fairly distant from Lepidoptera evolutionarily (3 mya), the genome analysis of species closely related to Lepidoptera has not yet been performed. The domesticated silkworm, *Bombyx mori*, has been used as a model for basic studies, which provides a number of mutants and genetically improved strains. In addition, several groups have engaged in the construction of molecular linkage maps in the silkworm using a variety of markers, with the aim of providing a framework for positional cloning of specific genes and mutations, large-scale physical map construction, analysis of quantitative trait loci (QTLs) and comparative genomics. To date, about 1,500 markers based on RAPDs, RFLPs, SSRs and ISSRs (Inter Simple Sequence Repeats) are available for the construction of molecular linkage maps, which now cover all 28 *Bombyx* chromosomes at an average spacing of 2 cM, equivalent to about 500 kb (refs 2–4). The robust genetic resources of *B. mori* make it an ideal reference for the Lepidoptera, where comparative genetics and genomics can work together to elucidate conserved evolutionary pathways and their diversification, identify new genes and gene systems as targets for transgenesis. This will also help basic research leading towards new genome-based approaches for the control of pest species. Therefore, the analysis of the *Bombyx* genome is one of the most urgent requirements in insect science today.

Aiming for the complete analysis of the *Bombyx* genome, we are taking the following strategies: (1) the construction of an Expressed Sequence Tags database (EST db), (2) the construction of a *Bombyx* BAC library, (3) making BAC contigs based on EST markers, and (4) genomic sequencing by BAC shot-gun sequencing. In this paper, we report on the progress of the cDNA project as the first step to the *Bombyx* genome analysis and the applications of the EST db.

## Strategy for preparing a comprehensive *Bombyx* cDNA catalog

A cDNA catalog is the comprehensive identification of all expressed genes by large-scale cDNA sequencing.

Unlike uni-cellular organisms such as yeast[5], the gene expression patterns in multi-cellular species deeply depend on tissues as well as developmental stages. The cDNAs from which ESTs are derived are present in libraries in proportion to the levels of mRNA in the tissue from which the library was prepared. Thus, ESTs are subject to 'expression bias' for multi-cellular species[6]. Therefore, we took the following strategy: We prepared many cDNA libraries of various tissues and different developmental stages, and then carried out random sequencing of a large number of cDNA clones. The use of a large number of cDNA libraries is effective to cover almost the whole set of *Bombyx* genes. In addition, this approach explicitly represents the tissue- and stage-specific gene expression patterns of all genes identified. Another advantage of this method is to represent all members of related genes and identify all members participating in the pathway of each biological process that the cells (or tissues) employ.

Figure 1 shows the flow chart of the *Bombyx* cDNA catalog. First, poly(A+)RNAs were extracted from many tissues, followed by cDNA synthesis with oligo-dT primers. To ensure the greatest effectiveness in gene classification by protein homology search, cDNA libraries were made by the directional cloning method (Lambda Zap

cDNA cloning kit, Stratagene). So far, 32 cDNA libraries of various tissues and different developmental stages were constructed (Table 1). More than 1,000 cDNA clones were picked up randomly and approximately 700 nucleotides of sequence from the 5′ end of the cDNA was determined. By a comparison of the deduced amino acid sequences with public protein databases such as PIR and Swissprotein databases, the gene classification was determined using a criterion of homology of more than 30% identity in a sequence greater than 100 amino acids. This is followed by a nucleotide homology search in the *Bombyx* EST database to give the expression profile of each gene.
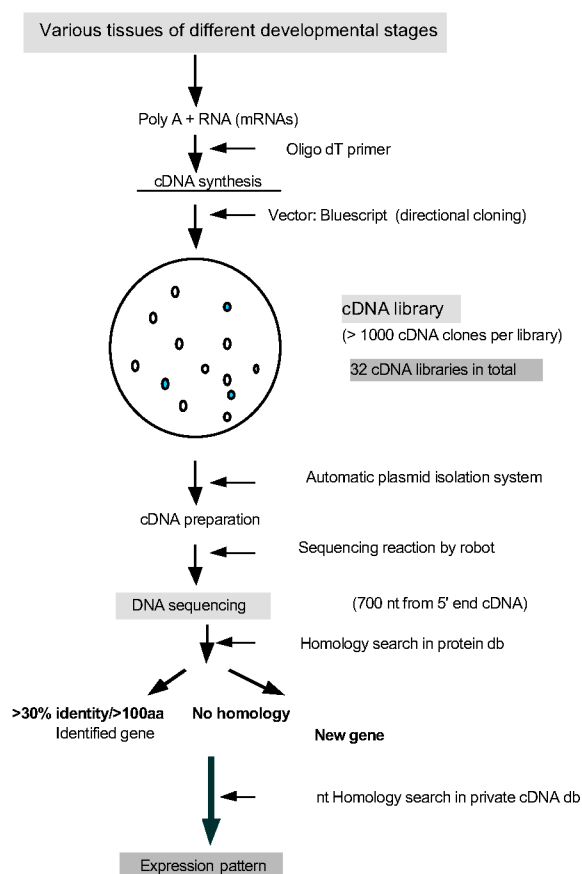
## Compilation of the EST database

We have sequenced about 26,000 cDNA clones from 32 cDNA libraries to date. Among those ESTs, we have identified about 9,500 non-redundant ESTs, which may cover more than 40% of all the genes of *Bombyx*, since
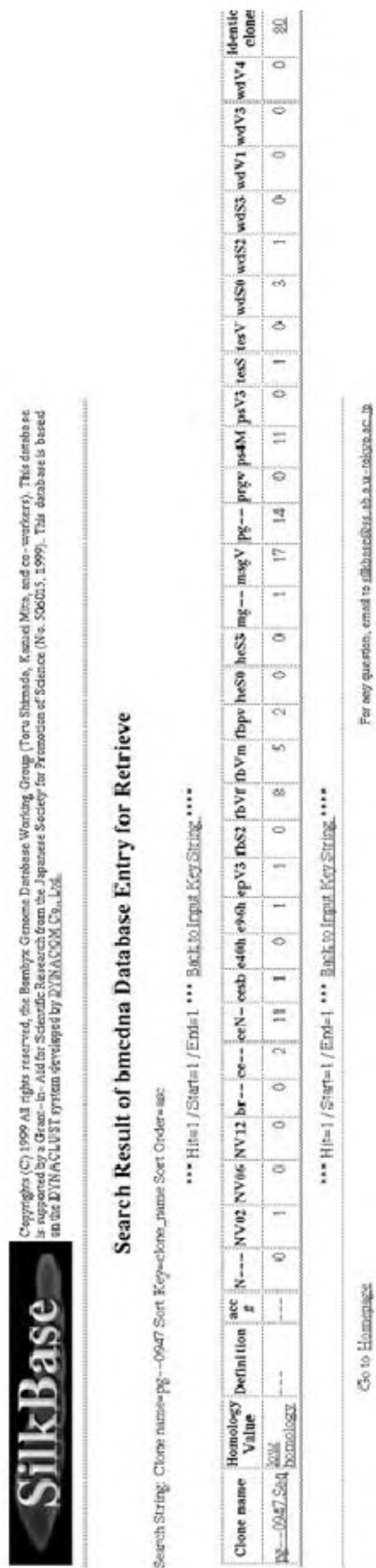
**Table 1.** List of *Bombyx* cDNA libraries for EST database

| Library name | Tissue/developmental stage |
|---|---|
| N--- | Culture cell BmN |
| NV02 | Baculovirus-infected BmN, 2 h post-infection |
| NV06 | Baculovirus-infected BmN, 6 h post-infection |
| NV12 | Baculovirus-infected BmN, 12 h post-infection |
| br— | Brain at the beginning of pupation |
| brS- | Brain, a mixture from day 3 of 5th larval instar to S2 |
| brP- | Brain, a mixture of pupal stages after P1 |
| ce— | Compound eyes, mixed stages from 5th instar larva to pupa |
| ceN- | Same as above but a different preparation |
| e40h | Embryo, 40 h after fertilization |
| e96h | Embryo, 96 h after fertilization |
| epV3 | Epidermis, 5th instar larva day 3 |
| fbS2 | Fat body, day 2 after spinning |
| fbVf | Female fat body, 5th instar larva day 3 |
| fbVm | Male fat body, 5th instar larva day 3 |
| fbpv | Baculovirus-infected fat body of S2 (2 days after spinning), 2 h post-infection |
| heS0 | Hemocytes, S0 (beginning of spinning) stage |
| heS3 | Hemocytes, S3 (3 days after spinning) stage |
| mg— | Midgut, 5th instar larva |
| msgV | Middle silkgland, 5th instar larva |
| pg— | Pheromone gland, adult moth |
| prgv | Prothoracic gland, 5th instar larva |
| ps4M | Posterior silkgland, 4th molt |
| psV3 | Posterior silkgland, 5th instar larva day 3 |
| tesS | Testis, spinning stage |
| tesV | Testis, 5th instar larva |
| wdV1 | Wing disc, 5th instar larva day 1 |
| wdV3 | Wing disc, 5th instar larva day 3 |
| wdV4 | Wing disc, 5th instar larva day 4 |
| wdS0 | Wind disc, the beginning of spinning |
| wdS2 | Wing disc, 2 days after spinning |
| wdS3 | Wing disc, 3 days after spinning |

Stages: P0, beginning of pupation; S0, beginning of spinning, S2, 2 days after spinning, S3, 3 days after spinning.



**Figure 1.** Scheme of construction of *B. mori* EST database.

**Figure 2.** *B. mori* EST database, 'SilkBase'. Explanation for each column is described in the text. Names of cDNA libraries are given in Table 1. 'SilkBase' is available at the website http://www.ab.a.u-tokyo.ac jp/silkbase/.

we roughly estimate that the total number of genes is 20,000. All sequenced ESTs are compiled into the EST database (Figure 2). The first column presents the clone name and the second column shows the results of the homology search. The third one gives the name of homologous gene, followed by its accession number in the fourth column. The subsequent columns present the numbers of identical clones found in each library, from which an information on the expression profile of each gene will be observed, i.e. a constitutive, tissue-specific or developmental stage-specific expression.

We have opened the whole *Bombyx* EST database named 'SilkBase' at the website http://www.ab.a.u-tokyo.ac.jp/silkbase/. 'SilkBase' is equipped with functions such as keyword/clone name search and BLAST search, which are very effective to search for homologous *Bombyx* cDNAs with known amino acid sequences of other species. Recently, another function for the comparison with FlyBase and WormBase for which genome sequencing has been completed was added to SilkBase. This will provide a powerful tool to screen candidates for Lepidoptera-specific genes.

## cDNA libraries

Table 1 presents a list of the cDNA libraries analysed. To study the changes of gene expression patterns during metamorphosis and development, several libraries were made from successive developmental stages of the same tissue, including brain, embryo, fat body, haemocyte, silkgland, testis and wing disc. For fat body, in addition to 5th instar larval (fbV) and pupal stages (fbS2), the sex difference in gene expression was observed by a comparison between libraries derived from 5th instar larval male (fbVm) and female fat body (fbVf). A comparison of the gene expression patterns between two imaginal discs, compound eye and wing disc, will reveal the factors responsible for tissue-specific gene expression caused by ecdysteroid hormone. For imaginal disc development, cDNA libraries of six successive stages of wing disc, wdV1, wdV3, wdV4, wdS0, wdS2 and wdS3, were prepared to analyse the genes related to tissue differentiation during metamorphosis[7,8]. For the compound eye cDNA library, not only is the tissue vanishingly small, but it also took so much time to dissect out that it was difficult to prepare enough tissue for each successive stage. Therefore, we made the cDNA library of mixed stages for the compound eye, and sequenced more than 5,000 cDNA clones. For a brain cDNA library, having a similar situation as the compound eye, we prepared three libraries, including a mixture of stages from day 3 of the 5th larval instar to day 3 post-spinning (brS), beginning pupation stage (P0), and pupal stage (brP). For the silk gland, comparison of a posterior silk gland cDNA libraries between 4th molt and the 5th instar will give the

genes responsible for molting, while the comparison between middle and posterior silkglands of fifth instar larvae may provide valuable information on tissue-specific gene expression. We also prepared cDNA libraries of baculovirus-infected culture cells at 0, 2, 6 and 12 h post-infection to understand the effects of virus infection on the gene expression of host cells and the mechanism of biological defense of host cells against virus infection.

## Characterization of cDNA libraries and approaches to increase coverage

Table 2 presents the average size of cDNA clones in each library. The size of the cDNA was estimated from agarose gel electrophoresis of PCR products amplified with vector primers. The average size of the cDNA was 0.92–1.59 kb, which indicated the high quality of the cDNA libraries analysed. Table 2 also shows the fraction of library-specific ESTs in the complete set of ESTs in each cDNA library, in which the libraries are arranged in the order of the time they were analysed. The values fluctuate between 25 and 54% showing no correlation with the time of analysis. This indicates that we have not yet reached saturation with the identical clones that have already been found in other cDNA libraries. However, it is obvious that saturation will have to become a serious problem as the ESTs accumulate more extensively. Therefore, we tried subtraction methods for compound eye (ce) and cultured cells (BmN). For the ce cDNA library, we dotted about 100 cDNA clones onto a nylon

**Table 2.** Characterization of cDNA libraries

| Library name | EST total number | Average size (kb) | Fraction of library-specific ESTs (%) |
|---|---|---|---|
| N--- | 756 | 1.10 | 27.4 |
| mg— | 604 | 0.92 | 39.6 |
| e40h | 756 | 1.58 | 39.4 |
| e96h | 704 | 1.39 | 35.5 |
| pg— | 468 | 1.38 | 34.0 |
| wdV4 | 1897 | 1.40 | 29.6 |
| wdS0 | 864 | 1.21 | 34.7 |
| wdS2 | 760 | 1.57 | 26.3 |
| wdS3 | 769 | 1.20 | 32.3 |
| prgv | 779 | 1.22 | 54.5 |
| msgV | 633 | n.d. | 40.3 |
| heS0 | 828 | 1.15 | 25.0 |
| heS3 | 588 | 0.95 | 30.6 |
| fbS2 | 281 | 1.44 | 38.4 |
| fbpv | 688 | 1.21 | 29.6 |
| ce— | 1578 | 1.02 | 36.8 |
| an— | 606 | 1.34 | 36.8 |
| brP- | 1587 | 1.40 | 33.1 |
| brS- | 798 | n.d. | 32.5 |

The cDNA libraries are arranged in the order of the time when the library was analysed, i.e. top, oldest; bottom, latest; n.d, not determined.

**Table 3.** Effects of normalization and subtraction on cDNA libraries

| Library name | EST total number | Library-specific EST (%) | Ribosomal proteins (%) | Tubulin (%) | Actin (%) | Elongation factors (%) |
|---|---|---|---|---|---|---|
| N--- | 752 | 27.4 | 9.4 | 1.1 | 1.5 | 2.8 |
| Nnor | 221 | 39.0 | 5.4 | 0.9 | 0.9 | 0 |
| ce— | 1578 | 36.8 | 3.6 | 3.9 | 12.4 | 0.6 |
| cesb | 83 | 43.4 | 3.1 | 2.4 | 1.3 | 0 |

Explanations for Nnor and cesb cDNA libraries are given in the text.

membrane filter (GeneScreen, NEN) and carried out hybridization using a mixed probe with 20 house-keeping genes which are expressed highly and ubiquitously such as several ribosomal protein genes, elongation factor 1 alpha gene, actin gene, etc. We picked up the clones which did not hybridize strongly with the mixed probe and sequenced them denoted cesb. In another approach, we employed a normalization method for cultured cell cDNA library denoted as Nnor[9,10]. A comparison between the original library and subtracted one using both methods is presented in Table 3. In both cases, the fraction of library-specific ESTs was increased, while the fractions of house-keeping genes examined were decreased by the subtraction procedures, indicating that these approaches efficiently increase the coverage of genes.

## Overview of results obtained from the silkworm EST database

Several interesting results such as identification of novel genes and comprehensive cloning of some gene families have already emerged from the 'SilkBase' search. In wing development, novel ecdysteroid-inducible genes were identified[11,12] and altogether we found nine cuticle protein genes expressed in the prepupal wing disc cDNA libraries[13]. A couple of important sex-determining factors including the doublesex gene were identified, which can provide clues to understanding the sex-determination mechanism in the silkworm and critical differences from *Drosophila*[14,15]. From the pheromone gland cDNA library we have identified homologs for acyl-CoA desaturase and acyl-CoA-binding protein which play a significant role in the production of the sex pheromones regulated by the neurohormone (PBAN)[16,17]. Comparative EST analysis of cDNA libraries derived from BmNPV-infected BmN cells revealed that the expression of several genes, including cytochrome C oxidase 1, increased in the late stages of virus infection, although most of the host genes were depressed as the infection progressed. In addition, two apoptosis-related genes of the host cells were identified during virus infection[18]. We also identified all members of the Bm tubulin alpha and beta gene families (Kawasaki *et al.*, submitted), and provided a list of novel proteinase inhibitors associated with *B. mori* cocoons[19].

Many other interesting genes have been identified in 'SilkBase', which have provided tools for homology-based identification of related genes in other insects, as well as phylogenetic analysis. These include a bacteria-induced serine proteinase inhibitor serpin gene found in *Manduca sexta*[20], and heat shock protein and related genes in *Spodoptera frugiperda*[21,22].

## Applications of the *Bombyx* EST database

### Construction of *B. mori* BAC contigs and physical map

We have constructed a high quality *B. mori* BAC library in collaboration with Pieter de Jong's group (Children's Hospital Oakland Research Institute, USA). Its average insertion size was estimated to be 168 kb with 11 fold redundancy. High molecular weight genomic DNA of *B. mori* (p50 strain) was extracted from isolated nuclei of posterior silkglands of day 2 of 5th instar larvae[23]. Construction of the *Bombyx* BAC library was exactly followed by the protocol of Osoegawa *et al.*[24]. We have begun to construct BAC contigs by filter hybridization using non-redundant ESTs as probes. The genome size of *B. mori* is estimated to be 530 Mb (refs 25, 26). If 6,200 non-redundant ESTs are available, one BAC clone will have 2–3 EST markers on average. Therefore, more than 6,200 non-redundant ESTs are needed to construct the complete set of BAC contigs. An important advantage of using this approach is that the physical map made by this method is collinear to the gene map since ESTs are used as markers. This offers the greatest validity that the ESTs can serve as anchors to the genetic map for positional cloning of mutations and for analysis of comparative genome organization of other Lepidoptera and other insects.

### EST microarray

Another major application of the EST database is EST microarrays. An EST microarray is a powerful tool for many functional studies as well as for genome analysis since it can provide quantitative expression profiles of a

large number of genes at one go. From the hybridization experiments using high density replica filters with EST probes to make BAC contigs, it was revealed that more than 10% of *B. mori* genes include repetitive sequences such as Bm1 in their 3′ UTR. Therefore, we designed and synthesized 6,000 specific primers located about 500 bp downstream from the 5′ end of cDNA to remove repetitive sequence from DNAs spotted on glasses. 6,000 DNAs were amplified by PCR with T3 vector primer and a specific primer, followed by spotting on glass slides. By this procedure, we could obtain almost equal sizes and amounts of DNA from PCR, resulting in even efficiency of DNA fixation on glass.

## Lepidoptera-specific genes

The Lepidopteran insects to which silkworm belongs have taxonomically specific biological phenomena including sex-determination, pheromone-dependent sexual communication, silk production, diapause, and insect–plant interactions, insect–microbe interactions, etc. Because the Lepidoptera include the most highly destructive agricultural pests, it is one of the most urgent targets to design novel insecticides today that should be effectively active to Lepidopteran pests only, but more harmless to other species and environment (see International Lepidopteran Genome Project Proposal, 2001). This should be effectively achieved by a better understanding through Lepidoptera-specific genes and their functions. Comparative studies on gene sequences among Lepidoptera, *Drosophila*, *C. elegans* and other species for which genome analysis is well advanced, can point to Lepidoptera-specific genes. An EST database in a model lepidopteran species will help to find gene sequences and gene functions not only in the model species itself, but also in non-model lepidopteras. Although no complete lepidopteran genome has yet been sequenced, the silkworm EST db with approximately 9,500 independent cDNAs will provide tools for analysis of silkworm genome. The comparison of silkworm ESTs with FlyBase and WormBase, that has been added to 'SilkBase' as a novel tool, is very valuable and useful in illustrating at the molecular level what makes lepidoptera different from other insects, and providing potential candidates for targets of lepidoptera-selective insecticides.

1. International Consortium for Lepidopteran Genomics, Lyon, 2001, August 16–17 (International Lepidopteran genome project proposal) http://www.ab.a.u-tokyo.ac.jp/lep-genome/.
2. Promboon, A., Shimada, T., Fujiwara, H. and Kobayashi, M., *Genet. Res.*, 1995, **66**, 1–7.
3. Yasukochi, Y., *Genetics*, 1998, **150**, 1513–1525.
4. Nagaraju, J., Klimenko, V. and Couble, P., in *Encyclopedia of Genetics* (ed. Reeve, E. C. R.), Fitzroy Press, London, 2001, pp. 219–239.
5. Morimyo, M. *et al.*, in *Biodefence Mechanisms Against Environmental Stresses* (eds Ozawa, T., Hori, T. and Tatsumi, K.), Kodansha Scientific, Tokyo, 1997, pp. 115–123.
6. Marra, M. A., Hillier, L. and Waterston, R. H., *Trends Genet.*, 1998, **14**, 4–7.
7. Chareyre, P., Guillet, C., Besson, M. T., Fourche, J. and Bosquet, G., *Insect Mol. Biol.*, 1993, **2**, 239–246.
8. Fletcher, J. C. and Thummel, C. S., *Development*, 1995, **121**, 1411–1421.
9. Soares, M. B., Bonaldo, M. F., Jelene, P., Su, L., Lawton, L. and Efstratiadis, A., *Proc. Natl. Acad. Sci. USA*, 1994, **91**, 9228–9232.
10. Bonaldo, M. F., Lenon, G. and Soares, M. B., *Genome Res.*, 1996, **6**, 791–806.
11. Quan, G-X., Mita, K., Okano, K., Shimada, T., Ugajin, N., Zhao, X., Goto, N., Kanke, E. and Kawasaki, H., *Insect. Biochem. Mol. Biol.*, 2000, **31**, 97–103.
12. Zhao, X., Mita, K., Shimada, T., Okano, K., Kanke, E. and Kawasaki, H., *ibid*, 2001, **31**, 1213–1219.
13. Takeda, M., Mita, K., Quan, G-X., Shimada, T., Okano, K., Kanke, E. and Kawasaki, H., *ibid*, 1019–1028.
14. Ohbayashi, F., Suzuki, M. G., Mita, K., Okano, K. and Shimada, T., *Comp. Biochem. Physiol.*, 2001, **B128**, 145–158.
15. Suzuki, M. G., Ohbayashi, F., Mita, K. and Shimada, T., *Insect. Biochem. Mol. Biol.*, 2001, **31**, 1201–1211.
16. Yoshiga, T., Okano, K., Mita, K., Shimada, T. and Matsumoto, S., *Gene*, 2000, **246**, 339–345.
17. Matsumoto, S. *et al.*, *Insect. Biochem. Mol. Biol.*, 2001, **31**, 603–609.
18. Okano, K., Shimada, T., Mita, K. and Maeda, S., *Virology*, 2001, **282**, 348–356.
19. Nirmala, X., Mita, K., Vanisree, V., Zurovec, M. and Sehnal, F., *Insect. Mol. Biol.*, 2001, **10** (in press).
20. Gan, H., Wang, Y., Jiang, H., Mita, K. and Kanost, M. R., *Insect. Biochem. Mol. Biol.*, 2001, **31**, 887–898.
21. Landais, I., Pommet, J-M., Mita, K., Nohata, J., Gimenez, S., Fournier, P., Devauchelle, G., Duonor-Cerutti, M. and Ogliastro, M., *Gene*, 2001, **271**, 348–356.
22. Lee, J., Hahn, Y., Yun, J. H., Mita, K. and Chung, J. H., *Biochim. Biophys. Acta*, 2000, **1491**, 355–363.
23. Ichimura, S., Mita, K., Zama, M. and Numata, M., *Insect. Biochem.*, 1985, **15**, 277–283.
24. Osoegawa, K., Woon, P. Y., Zhao, B., Frengen, E., Tateno, M., Catanese, J. J. and deJong, P. J., *Genomics*, 1998, **52**, 1–8.
25. Gage, L. P., *Chromosoma*, 1974, **45**, 27–42.
26. Rasch, E. M., in *Advances in Microscopy* (eds Cowden, R. R., Harrison, S. H. and Liss, A. R.), New York, 1985, pp. 137–166.